# 78. Information and Statistics. I

By Yukiyosi KAWADA

Department of Mathematics, Faculty of Science,
University of Tokyo

(Communicated by Shokichi IYANAGA, M. J. A., Sept. 14, 1987)

This short note is a summary of an axiomatic consideration of an information and its applications to statistics.[1]

**I. Axioms of an information 1.** Under an information we usually understand the *Kullback-Leibler information* (1951):

$$(1)\qquad I_{KL}(p, q) = \sum_{k=1}^{m} p_k \log(p_k/q_k)$$

where $p = (p_1, \cdots, p_m)$, $q = (q_1, \cdots, q_m)$ are two finite probability distributions. There are, however, similar known functions $I(p, q)$ which may be also called informations. For example,

$$(2)\qquad I_P(p, q) = \left(\sum_{k=1}^{m} p_k^2 q_k^{-1}\right) - 1 \quad \textit{(Pearson's information, 1900)}$$

which can be expressed as $(\sum_{k=1}^{m} (n_k - nq_k)^2 / nq_k)/n$ when $p = (n_1/n, \cdots, n_m/n)$ $(n = n_1 + \cdots + n_m)$, and

$$(3)\qquad I_K(p, q) = 2\left(1 - \sum_{k=1}^{m} p_k^{1/2} q_k^{1/2}\right) \quad \textit{(Kakutani's information, [5], 1948)}.$$

These are included as special cases of the family

$$(4)\qquad I^\lambda(p, q) = \frac{1}{\lambda}\left\{\left(\sum_{k=1}^{m} p_k^{1+\lambda} q_k^{-\lambda}\right) - 1\right\}, \quad -\frac{1}{2} \leq \lambda < \infty, \quad \lambda \neq 0,$$

namely, $I_P = I^1$, $I_K = I^{-1/2}$, and we define $I^0 = I_{KL}$.

We can easily see that

$$I^\lambda(p, q) \leq I^\mu(p, q) \qquad \text{for } \lambda < \mu$$

and

$$\lim_{n \to \infty} I^{\lambda_n}(p, q) = I^{\lambda_0}(p, q) \qquad \text{for } \lim_{n \to \infty} \lambda_n = \lambda_0.$$

We call $I^0$ the *parabolic* information, $I^\lambda$ $(\lambda > 0)$ a *hyperbolic* information and $I^{-\mu}(1/2 \geq \mu > 0)$ an *elliptic* information.

**Remark.** ( i ) The function $I^{-\mu}(p, q)$ for $1/2 \geq \mu > 0$ was introduced by several authors [4], [7] and the general case was also considered in [9].

( ii ) In the definition (4) we can extend the value $\lambda$ for $\lambda < -1/2$ formally. Then we have

$$I^{-\lambda}(p, q) = -\frac{\lambda - 1}{\lambda} I^{\lambda-1}(q, p), \quad \lambda > 1$$

$$I^{-1}(p, q) = 0$$

$$I^{-\mu}(p, q) = \frac{1 - \mu}{\mu} I^{\mu-1}(q, p), \quad 1/2 < \mu < 1.$$

---

1) The details will be published in the Proceedings of the Institute of Statistical Mathematics (Tōkei Sūri) in Japanese.

**Theorem 1.**  *The function* $I^\lambda(p, q) = I^\lambda(p_1, \cdots, p_m; q_1, \cdots, q_m)$ $(-1/2 \leq \lambda < \infty)$ $(p_k \geq 0,\ q_k \geq 0,\ p_1 + \cdots + p_m = q_1 + \cdots + q_m = 1,\ m = 1, 2, \cdots)$ *satisfies the following conditions:*

( I )  *Reducibility.*  *If* $p_m = q_m = 0$, *then*

$$I^\lambda(p_1, \cdots, p_m; q_1, \cdots, q_m) = I^\lambda(p_1, \cdots, p_{m-1}; q_1, \cdots, q_{m-1}).$$

(II)  *Symmetry.*

$$I^\lambda(p_1, \cdots, p_m; q_1, \cdots, q_m) = I^\lambda(p_{i_1}, \cdots, p_{i_m}; q_{i_1}, \cdots, q_{i_m})$$

*for any substitution* $(i_1, \cdots, i_m)$ *of* $(1, \cdots, m)$.

(III)  *Non-negativity.*

$$I^\lambda(p, q) \geq 0$$

*for any* $p, q$, *and the equality holds if and only if* $p = q$.

(IV)  *Convexity and* (V) *Invariance.*

$$I^\lambda(p_1 + p_2, p_3, \cdots, p_m; q_1 + q_2, q_3, \cdots, q_m)$$
$$\leq I^\lambda(p_1, p_2, \cdots, p_m; q_1, q_2, \cdots, q_m)$$

*holds in general and the equality holds if and only if* $q_1/p_1 = q_2/p_2$.

(VI)  *Additivity and pseudo-additivity.*  *Let* $p \otimes p'$, $q \otimes q'$ *be the direct product distributions.*  *If* $\lambda = 0$ *the additivity holds:*

$$I^0(p \otimes p', q \otimes q') = I^0(p, q) + I^0(p', q').$$

*In general the pseudo-additivity holds:*

$$I^\lambda(p \otimes p', q \otimes q') = I^\lambda(p, q) + I^\lambda(p', q') + \lambda I^\lambda(p, q) \cdot I^\lambda(p', q').$$

(VII)  *Continuity.*  *If* $\lim_{n \to \infty} p_n = p_0$, *and* $\lim_{n \to \infty} q_n = q_0$, *then*

$$\lim_{n \to \infty} I^\lambda(p_n, q_n) = I^\lambda(p_0, q_0).$$

(VIII)  *Relativity.*  *Let* $p^* = (p_{kj})$, $q^* = (q_{kj})$ $(k = 1, \cdots, m; j = 1, \cdots, r_k)$ *be probability distributions.*  *Put* $p_k = \sum_{j=1}^{r_k} p_{kj}$, $q_k = \sum_{j=1}^{r_k} q_{kj}$ *and* $p = (p_k)$, $q = (q_k)$.  *Define the conditional probability:* $p^{(k)} = (p_{kj}/p_k)$, $q^{(k)} = (q_{kj}/q_k)$ $(j = 1, \cdots, r_k)$ *for* $k = 1, \cdots, m$.  *Then*

$$I^\lambda(p^*, q^*) = I^\lambda(p, q) + \sum_{k=1}^{m} p_k^{1+\lambda} q_k^{-\lambda} I^\lambda(p^{(k)}, q^{(k)}).$$

*For* $\lambda = 0$ *these properties are proved in Kullback* [8].

**Remark.**  We see easily that (VIII) implies (VI), and (III) and (VIII) imply (IV) and (V).

**Theorem 2.**  *Let us fix a constant* $\lambda$ $(-1/2 \leq \lambda < \infty)$ *and assume that a function* $I(p_1, \cdots, p_m; q_1, \cdots, q_m)$ *satisfies* (I) *reducibility,* (II) *symmetry,* (III) *non-negativity,* (VII) *continuity, and* (VIII) *relativity.*  *Then*

$$I(p, q) = c I^\lambda(p, q)$$

*for some constant* $c > 0$.

**2.**  Now we shall consider about "information" which we define by the following system of axioms.

**Definition 1.**  Let $p = (p_1, \cdots, p_m)$, $q = (q_1, \cdots, q_m)$ $(\sum_{k=1}^{m} p_k = \sum_{k=1}^{m} q_k = 1)$ be finite probability distributions $(m = 1, 2, \cdots)$.  A function $I(p, q) = I(p_1, \cdots, p_m; q_1, \cdots, q_m)$ is called an *information* if the function $I$ satisfies the axioms: (I) *reducibility,* (II) *symmetry,* (III) *non-negativity,* (IV) *convexity, and* (V) *invariance.*

The functions $I^\lambda(p, q)$ $(-1/2 \leq \lambda < \infty)$ are informations in the above

sense.   We can also define an information in the form $I(p, q) = F(I_1(p, q),$ $\cdots, I_r(p, q))$ by using a suitable function $F(x_1, \cdots, x_r)$ from known informations $I_1, \cdots, I_r$.   For example, $I = aI_1 + bI_2$, $I = aI_1^2 + bI_2^2$, $a > 0$, $b > 0$, etc.   In particular,

$$( 5 )_1 \qquad\qquad \tilde{I}^\lambda(p, q) = \frac{1}{\lambda} \log (1 + \lambda I^\lambda(p, q)), \quad \lambda > 0,$$

$$( 5 )_2 \qquad\qquad \tilde{I}^{-\mu}(p, q) = \frac{-1}{\mu} \log (1 - \mu I^{-\mu}(p, q)), \quad 0 < \mu < \frac{1}{2}$$

are also informations which satisfy the additivity.

Remark.   ( i ) $\tilde{I}^{-\mu}$ was introduced by Kudō [7] (1953) and $\tilde{I}^\lambda$ and $\tilde{I}^{-\mu}$ were also considered by Rényi [10] (1961).

( ii )   The functions $d(p, q) = (\sum_{k=1}^m |p_k - q_k|)/\sqrt{2}$ and $D(p, q) = (\sum_{k=1}^m (p_k - q_k)^2)^{1/2}$ are not informations in the above sense.

3.   Definition 2.   A continuous information is called *fundamental* if $I(p, q)$ can be expressed as
$$( 6 ) \qquad\qquad I(p, q) = L(p_1, q_1) + \cdots + L(p_m, q_m)$$
by a continuous function $L(x, y)$ defined for $0 \leqq x \leqq 1$ and $0 \leqq y \leqq 1$.

For example, $I_{KL}, I_P, I_K, I^\lambda$ are fundamental, but $\tilde{I}^\lambda, \tilde{I}^{-\mu}$ are not fundamental.

Theorem 3.   *In order that a function $I(p, q)$ defined by (6) is an information it is necessary and sufficient that*
( i )   $L(0, 0) = 0$, $L(1, 1) = 0$,
( ii )   *if* $p_1/q_1 = p_2/q_2 = (p_1 + p_2)/(q_1 + q_2)$ $(p_1 + p_2 \leqq 1,\ q_1 + q_2 \leqq 1)$, *then*
$$L(p_1, q_1) + L(p_2, q_2) = L(p_1 + p_2,\ q_1 + q_2),$$
( iii )   *if* $p_1 + p_2 \leqq 1$, $q_1 + q_2 \leqq 1$, *then*
$$L(p_1, q_1) + L(p_2, q_2) \geqq L(p_1 + p_2,\ q_1 + q_2)$$
*and the equality holds if and only if* $p_1/q_1 = p_2/q_2$.

Theorem 4.   *A fundamental information $I(p, q)$ can be expressed in the form*
$$( 7 ) \qquad\qquad I(p, q) = \sum_{k=1}^m p_k K(q_k/p_k)$$
*by a non-negative strictly convex function $K(x)$ with $K(1) = 0$, and conversely the function defined by (7) is a fundamental information.   If we assume furthermore the differentiability of $K(x)$, such function $K(x)$ is uniquely determined by $I$.*

Examples.

$$I^0(p, q) = \sum_{k=1}^m p_k K^0(q_k/p_k), \qquad K^0(x) = -\log x + (x - 1) \geqq 0.$$

$$I^\lambda(p, q) = \sum_{k=1}^m p_k K^\lambda(q_k/p_k), \quad \lambda \neq 0, \quad K^\lambda(x) = (x^{-\lambda} - 1)/\lambda + (x - 1) \geqq 0.$$

Theorem 5.   *Let $I(p, q)$ be a differentiable fundamental information (i.e. $L(x, y)$ in (6) be three times differentiable in $x$ and $y$).*
( i )   *If $I(p, q)$ satisfies the additivity :*
$$I(p \otimes p',\ q \otimes q') = I(p, q) + I(p', q'),$$

*then we have*

$$I(p, q) = c_1 I^0(p, q) + c_2 I^0(q, p), \quad c_1 \geqq 0, \quad c_2 \geqq 0.$$

(ii)  *If $I(p, q)$ satisfies the relation*

$$I(p \otimes p', \, q \otimes q') = I(p, q) + I(p', q') + I(p, q) \cdot I(p', q')$$

*then we have*

$$I(p, q) = \lambda I^\lambda(p, q), \quad or \quad = \lambda I^\lambda(q, p)$$

*by some $\lambda > 0$.*

(iii)  *If $I(p, q)$ satisfies the relation*

$$I(p \otimes p', \, q \otimes q') = I(p, q) + I(p', q') - I(p, q) \cdot I(p', q')$$

*then we have*

$$I(p, q) = \mu I^{-\mu}(p, q), \quad or \quad = \mu I^{-\mu}(q, p)$$

*by some $\mu$ ($1/2 \geqq \mu > 0$).*

**Remark.**  Rényi [10] characterized $I^0$, $\tilde{I}^\lambda$ and $\tilde{I}^{-\mu}$ by different axioms.

Let $I$ be a differentiable fundamental information (7). Let $p = (p_k)$, $q = (q_k)$ and $q^0 = (q_k^0)$ be probability distributions and put $p_k = q_k^0 + u_k$, $q_k = q_k^0 + v_k$ ($k = 1, \cdots, m$).

If $|u_k| < \varepsilon$, $|v_k| < \varepsilon$ ($k = 1, \cdots, m$), then we have

(8)
$$I(p, q) = \frac{\alpha}{2} \sum_{k=1}^{m} \frac{1}{q_k^0} (u_k - v_k)^2 + R, \qquad R = O(\varepsilon^3),$$

where $\alpha = (d^2 K / dx^2)(1) \geqq 0$. We call $\alpha$ the *invariant* of $I$.

The invariant $\alpha$ of $I^\lambda$ is given by $\alpha(I^\lambda) = 1 + \lambda$ ($-1/2 \leqq \lambda < \infty$).

# References

[ 1 ]  H. Akaike:  Information theory and an extension of the maximum likelihood principle. 2nd Internat. Symp. on Information Theory, Akademiai Kiado, Budapest, pp. 267–281 (1973).

[ 2 ]  ——:  A new look at the statistical model identification. IFEE Trans. Automatic Control, AC-19, pp. 716–723 (1974).

[ 3 ]  S. M. Ali and S. D. Silvey:  A general class of divergence of one distribution from another. J. Roy. Statist. Soc., B **28**, 131–142 (1966).

[ 4 ]  I. Csizár:  Information measures; a critical survey. Trans. 7-th Prague Conference on Information Theory, Statistical Decision Functions, Random Processes and of the 1974 European Meeting of Statisticians, vol. B, pp. 73–86 (1978).

[ 5 ]  S. Kakutani:  On equivalence of infinite product measures. Ann. of Math., **49**, 214–224 (1948).

[ 6 ]  H. Kudō:  A theorem of Kakutani on infinite product measures. Nat. Sci. Rep. Ochanomizu Univ., **3**, 10–22 (1952).

[ 7 ]  ——:  Theory of time series and informations, and their applications, Chap. 2. Statistical Experiments and Their Informations. Nippon Kagaku-Gizyutu Renmei, pp. 104–124 (1953) (in Japanese).

[ 8 ]  S. Kullback:  Information Theory and Statistics, Wiley (1959).

[ 9 ]  P. N. Rathie and P. L. Kannappan:  A directed-divergence function of type $\beta^*$. Information and Control, **20**, 38–45 (1972).

[10]  A. Rényi:  On measures of entropy and information. Proc. Fourth Berkeley Symposium Math. Statist. and Probability, **1**, 547–561 (1961).