

86. Information and Statistics. II

By Yukiyoji KAWADA

Faculty of Science, University of Tokyo

(Communicated by Shokichi IYANAGA, M. J. A., Oct. 12, 1987)

This is a continuation of Kawada [0]. We use the same notations.

II. *L*-sets and informations. 1. Let $\mathbf{p}=(p_1, \dots, p_m)$ and $\mathbf{q}=(q_1, \dots, q_m)$ be probability distributions. We call the set

$$(9) \quad L(\mathbf{p}, \mathbf{q}) = \left\{ (x, y) \mid x = \sum_{k=1}^m \alpha_k p_k, y = \sum_{k=1}^m \alpha_k q_k, 0 \leq \alpha_k \leq 1, k=1, \dots, m \right\}$$

the *Liapunov-set* (simply *L-set*) of the pair (\mathbf{p}, \mathbf{q}) . See Kudō [6], [7].

$L(\mathbf{p}, \mathbf{q})$ has the following properties:

- (i) $L(\mathbf{p}, \mathbf{q}) = \Delta$ (the diagonal segment joining $(0, 0)$ and $(1, 1)$) if and only if $\mathbf{p} = \mathbf{q}$.
- (ii) $L(\mathbf{p}, \mathbf{q})$ contains the points $(0, 0)$ and $(1, 1)$.
- (iii) $L(\mathbf{p}, \mathbf{q})$ is contained in the square $[0, 1] \times [0, 1]$.
- (iv) $L(\mathbf{p}, \mathbf{q})$ is a symmetric convex set with the center $(1/2, 1/2)$.
- (v) Let the indices of (p_k, q_k) be so substituted that

$$0 \leq (q_1/p_1) \leq (q_2/p_2) \leq \dots \leq (q_m/p_m) \leq \infty$$

holds. Then

$$L(\mathbf{p}, \mathbf{q}) = \{(x, y) \mid \varphi(x) \leq y \leq \psi(x), 0 \leq x \leq 1\}$$

where $\varphi(x)$ is a polygon with $m+1$ vertices

$$(0, 0), (p_1, q_1), (p_1+p_2, q_1+q_2), \dots, (p_1+\dots+p_{m-1}, q_1+\dots+q_{m-1}), (1, 1)$$

and $\psi(x)$ is a polygon with $m+1$ vertices

$$(0, 0), (p_m, q_m), (p_m+p_{m-1}, q_m+q_{m-1}), \dots, (p_m+p_{m-1}+\dots+p_2, q_m+q_{m-1}+\dots+q_2), (1, 1).$$

Theorem 6. A function $I(\mathbf{p}, \mathbf{q})$ for any pair of finite probability distributions (\mathbf{p}, \mathbf{q}) is an information if and only if

- (i) $L(\mathbf{p}, \mathbf{q}) = \Delta \Rightarrow I(\mathbf{p}, \mathbf{q}) = 0$,
- (ii) $L(\mathbf{p}, \mathbf{q}) = L(\mathbf{p}', \mathbf{q}') \Rightarrow I(\mathbf{p}, \mathbf{q}) = I(\mathbf{p}', \mathbf{q}')$,
- (iii) $L(\mathbf{p}, \mathbf{q}) \supsetneq L(\mathbf{p}', \mathbf{q}') \Rightarrow I(\mathbf{p}, \mathbf{q}) > I(\mathbf{p}', \mathbf{q}')$.

Namely, an information I is characterized by the property that I is a monotone functional of the family of all *L*-sets with $I=0$ for $L=\Delta$.

2. (i) We can characterize a fundamental information I geometrically as

$$(10) \quad I_K(\mathbf{p}, \mathbf{q}) = \int_C K(d\varphi/dx) dx$$

where $K(x)$ is a non-negative differentiable function with $K(1)=K'(1)=0$, $K''(x)>0$, $\varphi(x)$ is the polygon defined as above and the integral is the curvilinear integral along the polygon $C: y=\varphi(x)$.

In particular, if we put

$$K(x) = \sqrt{1+x^2} - (x+1)/\sqrt{2},$$

then $I_K(\mathbf{p}, \mathbf{q}) = (\text{the length of the polygon } C) - \sqrt{2}$.

Thus we call the information (10) of the *type of arc-length*.

We can define several other types of informations geometrically.

(ii) Type of area of L -sets. Let

$$(11) \quad I_A(\mathbf{p}, \mathbf{q}) = \text{the area of } L\text{-set } L(\mathbf{p}, \mathbf{q}).$$

Then I_A is an information by Theorem 6. We can write also

$$I_A(\mathbf{p}, \mathbf{q}) = \left(\sum_{k=1}^m \sum_{l=1}^m |p_k q_l - p_l q_k| \right) / 2.$$

If we take a continuous function $f(x, y)$ defined on $0 \leq x \leq 1, 0 \leq y \leq 1$ and positive for $0 < x < 1, 0 < y < 1$, then

$$(12) \quad I_{A,f}(\mathbf{p}, \mathbf{q}) = \iint_{L(\mathbf{p}, \mathbf{q})} f(x, y) dx dy$$

is also an information by Theorem 6. We call these informations $I_{A,f}$ of the *type of area*.

(iii) Type of breadth of L -sets. Let $B_\theta(L)$ be the breadth of the convex set $L(\mathbf{p}, \mathbf{q})$ in the direction θ with the x -axis, and $f(\theta)$ ($0 \leq \theta < \pi$) be any positive continuous function. Then

$$(13) \quad I_{B,f}(\mathbf{p}, \mathbf{q}) = \frac{1}{\pi} \int_0^\pi (B_\theta(L) - B_\theta(\Delta)) f(\theta) d\theta$$

is an information by Theorem 6. We call these informations $I_{B,f}$ of the *type of breadth*. Notice that $d(\mathbf{p}, \mathbf{q}) = B_{3\pi/4}(L)$.

3. Now we introduce the concept of completeness of a family of informations after Kudō [7].

Definition 3. A family of informations $\{I_\omega(\mathbf{p}, \mathbf{q}) \mid \omega \in \Omega\}$ is called *weakly complete* if $I_\omega(\mathbf{p}, \mathbf{q}) = I_\omega(\mathbf{p}', \mathbf{q}')$ for all $\omega \in \Omega$ implies $L(\mathbf{p}, \mathbf{q}) = L(\mathbf{p}', \mathbf{q}')$, and is called *strongly complete* if $I_\omega(\mathbf{p}, \mathbf{q}) \geq I_\omega(\mathbf{p}', \mathbf{q}')$ for all $\omega \in \Omega$ implies $L(\mathbf{p}, \mathbf{q}) \supset L(\mathbf{p}', \mathbf{q}')$.

Any strongly complete family is evidently weakly complete, but the converse does not hold in general as shown by a counter example by K. Iseki in [7].

Theorem 7. (i) *The family of informations of the type of area :*

$$\left\{ I_A^{(i,j)}(\mathbf{p}, \mathbf{q}) = \iint_{L(\mathbf{p}, \mathbf{q})} x^i y^j dx dy \mid i, j = 0, 1, 2, \dots \right\}$$

is weakly complete, but not strongly complete.

(ii) *The family of informations of the type of area*

$$\left\{ I_{A,f}(\mathbf{p}, \mathbf{q}) = \iint_{L(\mathbf{p}, \mathbf{q})} f(x, y) dx dy \mid f(x, y) \geq 0 \text{ and continuous} \right\}$$

is strongly complete.

Theorem 8. (i) *The family of informations of the type of arc-length (i.e. fundamental informations)*

$$\{I^1(\mathbf{p}, \mathbf{q}) \mid \alpha < \lambda < \alpha + \varepsilon\} \quad (\alpha > 0, \varepsilon > 0)$$

and

$$\{I^{-\alpha}(\mathbf{p}, \mathbf{q}) \mid \alpha - \varepsilon < \mu < \alpha\} \quad (1/2 \geq \alpha > \varepsilon > 0)$$

are both weakly complete, but not strongly complete.

(ii) The family of informations of the type of arc-length

$$\left\{ I_{\kappa}(\mathbf{p}, \mathbf{q}) = \sum_{k=1}^m p_k K(q_k/p_k) \mid K(1) = K'(1) = 0, K''(x) > 0 \text{ for } x > 0 \right\}$$

is strongly complete.

III. Applications to statistics. 1. H. Akaike [1], [2] established the theory of AIC (Akaike information criterion), whose direct application gives a method of model selection from the standpoint of prediction. There he used as a basic tool the Kullback-Leibler information I_{KL} . Here we shall show that a similar results can be obtained if we use a regular information I , which we shall define below, instead of I_{KL} .

Definition 4. An information I is called *regular* if the following two conditions (A) and (B) hold.

(A) Let $\mathbf{p} = (p_1, \dots, p_m)$, $\mathbf{q} = (q_1, \dots, q_m)$ and $\mathbf{q}^0 = (q_1^0, \dots, q_m^0)$ be finite probability distributions, and put

$$p_k = q_k^0 + u_k, \quad q_k = q_k^0 + v_k \quad (k=1, \dots, m)$$

$$u_1 + \dots + u_m = 0, \quad v_1 + \dots + v_m = 0.$$

For $|u_k| < \varepsilon$, $|v_k| < \varepsilon$ ($k=1, \dots, m$) $I(\mathbf{p}, \mathbf{q})$ is three times differentiable with respect to $(u_1, \dots, u_{m-1}, v_1, \dots, v_{m-1})$ and

$$(14) \quad I(\mathbf{p}, \mathbf{q}) = \frac{\alpha}{2} \sum_{k=1}^m \frac{1}{q_k^0} (u_k - v_k)^2 + R, \quad R = O(\varepsilon^3)$$

holds, where α is a positive constant. α is called the *invariant* of I .

(B) If we fix \mathbf{q} then for any \mathbf{p} the inequality

$$0 \leq I(\mathbf{p}, \mathbf{q}) \leq c(\mathbf{q})$$

holds, where $c(\mathbf{q})$ is a certain constant. By (8) any differentiable fundamental information satisfies the condition (A), and we can easily verify that $I'(-1/2) < \lambda < \infty$ satisfies the condition (B).

Let $\mathbf{q}^0 = (q_1^0, \dots, q_m^0)$ be a probability distribution on m events (E_1, \dots, E_m) . Suppose that the events E_1, \dots, E_m occur N_1, \dots, N_m times respectively in n ($n = N_1 + \dots + N_m$) independent trials, and put

$$(15) \quad \mathbf{P} = (N_1/n, \dots, N_m/n).$$

Theorem 9. Let I be a regular information with the invariant α . Then as $n \rightarrow \infty$ the random variable $(2n/\alpha)I(\mathbf{P}, \mathbf{q}^0)$ converges in distribution to the chi-square distribution χ_{m-1}^2 with $m-1$ degrees of freedom. Moreover,

$$\lim_{n \rightarrow \infty} (2n/\alpha)E(I(\mathbf{P}, \mathbf{q}^0)) = m-1$$

holds, where E means the expectation of the random variable.

2. Now suppose that we are given a family of distributions $\mathbf{q}(\theta) = (q_1(\theta), \dots, q_m(\theta))$, $\theta = (\theta_1, \dots, \theta_r)$ ($\theta \in \Omega^{(r)}$) with r continuous parameters. We assume that the unknown true probability distribution \mathbf{q}^0 is contained in this family as $\mathbf{q}^0 = \mathbf{q}(\theta^0)$, $\theta^0 = (\theta_1^0, \dots, \theta_r^0)$. We define the random probability distribution \mathbf{P} by (15) after n independent trials. For this value \mathbf{P} , choose the value of parameters $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_r)$ such that $I(\mathbf{P}, \mathbf{q}(\hat{\theta}))$ takes its minimum at $\theta = \hat{\theta}$. We can consider $\hat{\theta}$ also as a random vector.

Theorem 10. *As $n \rightarrow \infty$ the random vector.*

$$\sqrt{n}(\hat{\theta}_1 - \theta_1^0, \dots, \hat{\theta}_r - \theta_r^0)$$

converges in distribution to the r -dimensional normal distribution $N((0, \dots, 0), \alpha I^{-1})$ with the mean vector $(0, \dots, 0)$ and the variance matrix αI^{-1} , where

$$I = \left(\left(\frac{\partial^2 I}{\partial \theta_i \partial \theta_j} \right)_{\theta = \theta^0} \right)_{i, j = 1, \dots, r} = \alpha Q \cdot Q,$$

$$Q = (q_j^{(i)} / \sqrt{q_j})_{\theta = \theta^0}, \quad q_j^{(i)} = \partial q_j / \partial \theta_i.$$

Theorem 11. *As $n \rightarrow \infty$ the random variable $(2n/\alpha)(I(\mathbf{P}, \mathbf{q}^0) - I(\mathbf{P}, \mathbf{q}(\hat{\theta})))$ converges in distribution to the chi-square distribution χ_r^2 with r degrees of freedom, and $(2n/\alpha)I(\mathbf{P}, \mathbf{q}(\hat{\theta}))$ itself converges in distribution to the chi-square distribution χ_{m-1-r}^2 with $m-1-r$ degrees of freedom.*

3. Now let $\mathbf{q}^0 = (q_1^0, \dots, q_m^0)$ be unknown true probability distribution of the events (E_1, \dots, E_m) , and suppose that we obtain the events $E_1, \dots, E_m, n_1, \dots, n_m$ times respectively in n ($n = n_1 + \dots + n_m$) independent trials. Put $\mathbf{p}^0 = (n_1/n, \dots, n_m/n)$.

Let $\Omega^{(r)} = \{\mathbf{q}(\theta) = (q_1(\theta), \dots, q_m(\theta))\}$ be a model for \mathbf{q}^0 which contains $\mathbf{q}^0 = \mathbf{q}(\theta^0)$. Assume that I is a regular information, and $\hat{\theta}$ is the value of θ in a neighbourhood of θ^0 such that $I(\mathbf{p}^0, \mathbf{q}(\hat{\theta}))$ is the minimum. Now define (16)

$$AIC(\Omega^{(r)}) = (2n/\alpha)I(\mathbf{p}^0, \mathbf{q}(\hat{\theta})) + 2r$$

after Akaike [1], [2]. Akaike's method of selection of model is as follows.

Suppose we are given several models for \mathbf{q}^0 . i.e. $\Omega^{(r_1)}, \dots, \Omega^{(r_s)}$. After n independent trials we obtain \mathbf{p}^0 as above. Compare the values $AIC(\Omega^{(r_t)})$ ($t = 1, \dots, s$). Choose the model $\Omega^{(r_t)}$ for which $AIC(\Omega^{(r_t)})$ takes the minimum among s values.

This method depends on the following theorem in prediction theory. Namely, we repeat n^* new independent trials, for which the events E_1, \dots, E_m occur N_1^*, \dots, N_m^* times ($n^* = N_1^* + \dots + N_m^*$) respectively. Put

$$\mathbf{P}^* = (N_1^*/n^*, \dots, N_m^*/n^*).$$

The mean value $E^*(I(\mathbf{P}^*, \mathbf{q}(\hat{\theta})))$ may be called the mean information in prediction.

Theorem 12.

$$AIC(\Omega^{(r)}) = (2n/\alpha)E^*(I(\mathbf{P}^*, \mathbf{q}(\hat{\theta}))) + R_1 + R_2$$

where

$$R_1 = (2n/\alpha)(I(\mathbf{p}^0, \mathbf{q}^0) - E^*(I(\mathbf{P}^*, \mathbf{q}^0)))$$

depends only on the value \mathbf{p}^0 , and R_2 is a random variable with $E(R_2) = 0$.

Reference^{*)}

- [0] Y. Kawada: Information and statistics. I. Proc. Japan Acad., 63A, 281-284 (1987).

^{*)} Other references are given in [0].