

Igor Vajda

On the statistical decision problems with finite parameter space

*Kybernetika*, Vol. 3 (1967), No. 5, (451)--466

Persistent URL: <http://dml.cz/dmlcz/124363>

## Terms of use:

© Institute of Information Theory and Automation AS CR, 1967

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these

*Terms of use.*



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library*  
<http://project.dml.cz>

## On the Statistical Decision Problems with Finite Parameter Space

IGOR VAJDA

In the paper there are obtained some both-sides estimates of the Bayes risk and information in a sample concerning a parameter in terms of total variations of pairs of conditional distributions in the framework of the classical model of statistical decision with finite parameter space. On the base of these relations the rate of convergence of the risk and of information to their limit values is studied, for the case if the sample size tends to infinity.

### 1. INTRODUCTION AND PRELIMINARIES

In this paper we shall deal with the classical model of statistical decision with a finite parameter probability space  $(X, \mathcal{X}, \mu)$  ( $\mathcal{X}$  is assumed to be the class of all subsets of the set  $X$  and the prior probability  $\mu(x)$  is assumed to be positive for every parameter value  $x \in X$ ), with measurable sample space  $(Y, \mathcal{Y})$ , set  $\{v_x\}$ ,  $x \in X$ , of conditional distributions of the variable  $y \in Y$  defined on the  $\sigma$ -algebra  $\mathcal{Y}$ , decision measurable space  $(X, \mathcal{X})$ , and with a non-negative loss function  $w$  defined on  $X \otimes X$ .

It is clear that a most important numerical characteristics of the model just described is the so called Bayes risk, i.e. the quantity defined as

$$(1.1) \quad r = \inf \sum_{x \in X} \mu(x) \int_Y w(x, \varrho(y)) dv_x(y),$$

where the infimum is extended over the set of all  $\mathcal{Y}$ -measurable decision functions  $\varrho : Y \rightarrow X$  (for a more particular discussion of the decision model, see [3]–[6]). Another characteristics of great importance is the average amount of information  $I$  contained in the sample  $y$  concerning the parameter  $x$  defined as

$$(1.2) \quad I = \int_{X \otimes Y} \log f(x, y) d\omega(x, y),$$

452 where  $f$  is the Radon-Nikodym density of the joint probability distribution

$$\omega(E) = \sum_{x \in X} \mu(x) v_x(\{y \in Y : (x, y) \in E\}), \quad E \in \mathcal{X} \otimes \mathcal{Y},$$

with respect to the Cartesian product distribution  $\mu \otimes \tilde{\omega} \gg \omega$  on  $\mathcal{X} \otimes \mathcal{Y}$ , where by  $\tilde{\omega}$  we denote the marginal distribution of  $\omega$  on  $\mathcal{Y}$ . (Let us note that all logarithms in this paper are taken to the base  $e$ .)

From the intuitive point of view it is clear that, though the information  $I$  does not depend on the loss function  $w$ , for a sufficiently wide class of loss functions there exists a relation between  $r$  and  $I$ . The data reduction theory developed by A. Perez [3], [4], shows that the indicated relation plays a growing role in solutions of certain class of decision problems.

The purpose of this paper is to estimate the Bayes risk and information in the framework of the decision model with finite parameter space and to investigate the indicated relation between them. This general questions are studied in Sec. 2; the results of Sec. 2 are then used in Sec. 3 devoted to the study of the rate of convergence of information and Bayes risk in some classes of decision models as it is more precisely described below.

Let the measurable sample space of the model we have considered be of the form

$$(1.3) \quad (Y^n, \mathcal{Y}^n) = \bigotimes_{i=1}^n (Y_i, \mathcal{Y}_i), \quad n = 1, 2, \dots, \infty,$$

(i.e. suppose that the samples are of size  $n$ ,  $y = (y_1, y_2, \dots, y_n)$ ), where  $Y_i$  is the set of all  $y_i$ 's with a given  $\sigma$ -algebra  $\mathcal{Y}_i$ , and denote the joint probability distribution of the  $n$ -vector  $(y_1, \dots, y_n)$  under the condition that  $x \in X$  is the realized value of the parameter by  $v_x^n$ . It is clear that, for every  $x \in X$ ,  $v_x^n$  is the restriction of  $v_x^\infty$  on the  $\sigma$ -algebra  $\mathcal{Y}^n \subset \mathcal{Y}^\infty$ , where the latter inclusion (as well as the inclusions  $\mathcal{Y}_i \subset \mathcal{Y}^n$ ,  $i = 1, 2, \dots, n$ ,  $n = 1, 2, \dots, \infty$ , that will be used below) is written in accordance with a well-known identification convention for product  $\sigma$ -algebras.

We shall denote by  $I_n$  or  $r_n$  the information or the Bayes risk respectively corresponding to the measurable space (1.3). It is known that  $I_n$  or  $r_n$ ,  $n = 1, 2, \dots$ , is a non-decreasing or non-increasing sequence respectively and that  $\lim I_n = I_\infty$ ,  $\lim r_n = r_\infty$  as  $n$  tends to infinity (cf. [3], [4]). From the point of view of applications it is important to ask which is the rate of convergence above. This question is studied in Sec. 3 under the assumption that, for every realized value  $x \in X$ , the sequence  $y_1, y_2, \dots$ , of samples is independent random sequence, i.e. that

$$(1.4) \quad v_x^\infty = \bigotimes_{i=1}^\infty v_{x_i} \quad \text{for every } x \in X,$$

where  $v_{x_i}$ ,  $i = 1, 2, \dots, \infty$ , denotes over all the paper the restriction of  $v_x^\infty$  on the sub- $\sigma$ -algebra  $\mathcal{Y}_i \subset \mathcal{Y}^\infty$ , and that it is a  $*$ -mixing random sequence, i.e. that there exists a positive integer  $N$  and a non-increasing real valued function  $\varphi$  defined for

all integers  $k \geq N$  with  $\lim \varphi(k) = 0$  for  $k \rightarrow \infty$  such that, for every  $k \geq N, l \leq m, m + k \leq n$  and  $x \in X$ , the following inequality takes place

$$(1.5) \quad |v_x^n(E \cap F) - v_x^n(E) v_x^n(F)| \leq \varphi(k) v_x^n(E) v(F)$$

for every

$$E = \{(y_1, \dots, y_n) : (y_l, y_{l+1}, \dots, y_m) \in \tilde{E}\}$$

and

$$F = \{(y_1, \dots, y_n) : (y_{m+k}, y_{m+k+1}, \dots, y_n) \in \tilde{F}\},$$

where

$$\tilde{E} \in \bigotimes_{i=l}^m \mathcal{A}_i$$

and

$$\tilde{F} \in \bigotimes_{i=m+k}^n \mathcal{A}_i.$$

We shall show that in both these cases the exponential rate of convergence takes place.

In the remainder of this section we list some properties of a  $\Delta$ -divergence of two probability distributions for references later. The concept of  $\Delta$ -divergence plays an important role over all this paper.

If  $\nu_1, \nu_2$  are probability measures on a measurable space  $(Y, \mathcal{A})$ , then  $\Delta$ -divergence  $\Delta(\nu_1, \nu_2)$  is defined by

$$(1.6) \quad \Delta(\nu_1, \nu_2) = \frac{1}{2} \int_Y |\xi_1 - \xi_2| d(\nu_1 + \nu_2),$$

where  $\xi_i$  is the  $\mathcal{A}$ -measurable version of the Radon-Nikodym density  $d\nu_i/d(\nu_1 + \nu_2)$  for  $i = 1, 2$ .  $\Delta$  is clearly reflexive, symmetric, and the triangle inequality satisfying distance measure in the space of all probability distributions on  $(Y, \mathcal{A})$ . It is clear that  $2\Delta(\nu_1, \nu_2) = |v_1 - v_2|(Y)$ , where  $|v_1 - v_2|$  denotes the total variation of the signed measure  $\nu_1 - \nu_2$ . This implies in particular that  $\Delta$  takes values between 0 and 1, and that  $\Delta(\nu_1, \nu_2) = 0$  if and only if  $\nu_1 = \nu_2$  (on  $\mathcal{A}$ ) and  $\Delta(\nu_1, \nu_2) = 1$  if and only if  $\nu_1 \perp \nu_2$ , and also the following very useful statement:

(i) There exists a set  $F \in \mathcal{A}$  such that

$$\nu_1(E) - \nu_2(E) \leq \nu_1(F) - \nu_2(F) = \Delta(\nu_1, \nu_2) \quad \text{for every } E \in \mathcal{A}.$$

The following two assertions follow from the theory of semi-martingales ([2], Th. 4.1s, Chap. VII), as  $|\xi|$  is convex function of  $\xi$  (cf. the definition of  $\Delta$ -divergence).

(ii) If  $\mathcal{A}^{(n)}, n = 1, 2, \dots$ , is a non-decreasing sequence of sub- $\sigma$ -algebras of the  $\sigma$ -algebra  $\mathcal{A}$  such that  $\mathcal{A}$  is generated by

$$\bigcup_{n=1}^{\infty} \mathcal{A}^{(n)}$$

and if we denote by  $v_1^{(n)}, v_2^{(n)}$  the restriction of the distributions  $v_1, v_2$  on  $\mathscr{Y}^{(n)}$ , then  $A(v_1^{(n)}, v_2^{(n)})$  is for  $n = 1, 2, \dots$  a non-decreasing sequence and

$$\lim_{n \rightarrow \infty} A(v_1^{(n)}, v_2^{(n)}) = A(v_1, v_2).$$

(iii) If  $(Z, \mathscr{Z})$  is a measurable space and if  $T$  is a measurable transformation of  $(Y, \mathscr{Y})$  to  $(Z, \mathscr{Z})$ , then

$$A(v_1 T^{-1}, v_2 T^{-1}) \leq A(v_1, v_2).$$

(iv) If there exist two numbers  $a, b \geq 0$  and probability measures  $v_{ij}, i, j = 1, 2$ , such that

$$v_1 = av_{11} + bv_{12},$$

$$v_2 = av_{21} + bv_{22},$$

where  $v_{i1} \perp v_{j2}, i, j = 1, 2$ , then

$$A(v_1, v_2) = a A(v_{11}, v_{21}) + b A(v_{12}, v_{22}).$$

**Proof.** The proof of this equality can be based on (i) and on the assumption of singularity. Details are omitted here.

(v) If  $B(\cdot; p, n), B(\cdot; q, n)$  are binomial distributions with  $p \neq q$ , then there exist numbers  $A > 0, 0 < \lambda < 1$  such that

$$A(B(\cdot; p, n), B(\cdot; q, n)) > 1 - A\lambda^n \quad \text{for every } n = 1, 2, \dots$$

**Proof.** According to (i), if  $p < q$ , then

$$A(B(\cdot; p, n), B(\cdot; q, n)) = 1 - \min_{k=0,1,\dots,n} (b_1(n, k) + b_2(n, k)),$$

where

$$b_1(n, k) = \sum_{i=k+1}^n \binom{n}{i} p^i (1-p)^{n-i}$$

$$b_2(n, k) = \sum_{i=0}^k \binom{n}{i} q^i (1-q)^{n-i}.$$

If one of the numbers  $p, q$  is equal 0 or 1, (v) is trivial. Let, for  $0 < p < q < 1$ ,  $k_n$  be the least integer greater than or equal to  $\varrho n$ , where  $0 < \varrho < 1$  is the unique solution of the equation

$$\left(\frac{p}{q}\right)^{\varrho} = \left(\frac{1-q}{1-p}\right)^{1-\varrho}.$$

Using Stirling's formula and some elementary properties of the numbers  $B(\cdot; \cdot, n)$  it can be shown that there is  $A > 0$  and  $0 < \lambda < 1$  such that for every  $n = 1, 2, \dots$

$$b_1(n, k_n) + b_2(n, k_n) < A\lambda^n;$$

the remainder of the proof is now clear.

In what follows we shall denote by  $\text{card}(X)$  the cardinal of  $X$  and by  $H(\mu)$ , in accordance with [6], the entropy of the finite probability space  $(X, \mu)$ .

## 2. GENERAL INEQUALITIES

**Lemma 1.** *There exists a disjoint system of sets  $\{E_x\}$ ,  $E_x \in \mathcal{O}$ ,  $x \in X$ , such that*

$$(2.1) \quad v_x(Y - E_x) < \text{card}(X) (1 - \min_{x' \neq x} A(v_x, v_{x'})),$$

$$(2.2) \quad v_x(E_{x'}) < 1 - A(v_x, v_{x'}) \quad \text{for every } x \neq x'.$$

*Proof.* According to (i), there exists for every  $x, x' \in X$  a set  $E_{xx'}$  such that

$$(2.3) \quad A(v_x, v_{x'}) = v_x(E_{xx'}) - v_{x'}(E_{xx'});$$

hence

$$A(v_x, v_{x'}) \leq v_x(E_{xx'})$$

or

$$(2.4) \quad v_x(Y - E_{xx'}) \leq 1 - A(v_x, v_{x'})$$

and consequently

$$v_x\left(\bigcup_{x' \neq x} (Y - E_{xx'})\right) < \text{card}(X) (1 - \min_{x' \neq x} A(v_x, v_{x'})).$$

Since

$$\bigcup_{x' \neq x} (Y - E_{xx'}) = Y - \bigcap_{x' \neq x} E_{xx'},$$

it remains to put in the latter inequality

$$(2.5) \quad E_x = \bigcap_{x' \neq x} E_{xx'}$$

and (2.1) is proved. From (2.3) we can easily obtain the following relation

$$(2.6) \quad E_{xx'} = Y - E_{x'tx} \quad \text{for all } x, x' \in X.$$

Since

$$v_x(E_{x'}) < v_x(E_{x'x}) = v_x(Y - E_{xx'}),$$

456 to prove (2.2) it remains to use (2.4). In order to prove that the system of measurable sets  $\{E_x\}$ ,  $x \in X$ , defined by (2.5) is disjoint let us point out the following inclusions

$$E_x = \bigcap_{x' \neq x} E_{xx'} \subset E_{xx'}$$

$$E_{x'} = \bigcap_{x'' \neq x'} E_{x'x''} \subset E_{x'x}$$

and then let us use (2.6).

**Theorem 1.** *If we denote*

$$(2.7) \quad A = \sqrt{\text{card}(X) \left(1 + 2 \sum_x \sqrt{[\mu(x)(1 - \mu(x))]} \right)},$$

$$(2.8) \quad \mu = \min_{x \in X} \mu(x) > 0,$$

then

$$(2.9) \quad \mu \log 2 \left(1 - \min_{x' \neq x} A(v_x, v_{x'})\right) \leq H(\mu) - I \leq A \sqrt{1 - \min_{x' \neq x} A(v_{xx}, v_{x'})}.$$

Proof. Since, for every  $x, x' \in X$  the following inequality holds

$$\frac{\mu}{2} \leq \frac{\mu(x)\mu(x')}{\mu(x) + \mu(x')}$$

the left inequality in (2.9) follows from Th. 3 in [6]. We next prove the right inequality. It follows from Lemma 2 in [6] that, for the class  $\{E_x\}$ ,  $E_x \in \mathcal{A}$ ,  $x \in X$ , defined in Lemma 1 above,

$$H(\mu) - I \leq \sum_{x, x' \in X} \sqrt{[\mu(x) v_x(E_{x'})] \sum_{x'' \neq x} (\mu(x'') v_{x''}(E_{x'}))} + \\ + \sum_{x \in X} \sqrt{[\mu(x) v_x(E_0)] \sum_{x'' \neq x} (\mu(x'') v_{x''}(E_0))},$$

where

$$E_0 = \bigcap_{x \in X} (Y - E_x),$$

and consequently we can write

$$H(\mu) - I \leq \sum_{x \in X} \sqrt{[\mu(x) v_x(E_x)] \sum_{x'' \neq x} \mu(x'') v_{x''}(E_x)} + \\ + \sum_{x \in X} \sum_{x' \neq x} \sqrt{[\mu(x) v_x(E_{x'})] \sum_{x'' \neq x} \mu(x'') v_{x''}(E_{x'})} + \\ + \sum_{x \in X} \sqrt{[\mu(x) v_x(Y - E_x)] \sum_{x'' \neq x} \mu(x'')}. \quad \square$$

If we denote the terms on the right side subsequently by (I), (II), and (III), then by (2.1)

$$(2.10) \quad (III) \leq \sqrt{[\text{card}(X) (1 - \min_{x \neq x'} d(v_x, v_{x'}))] \sum_{x \in X} \sqrt{[\mu(x) (1 - \mu(x))]}.$$

If we apply on the sum (I) the Schwarz's inequality, we obtain

$$(2.11) \quad (I) \leq \sqrt{\sum_{x' \in X} \mu(x') v_{x'}(Y - E_{x'})} \leq \sqrt{[\text{card}(X) (1 - \min_{x' \neq x} d(v_x, v_{x'}))]}.$$

One more application of Schwarz's inequality yields

$$\begin{aligned} (II) &\leq \sum_{x \in X} \sqrt{[\mu(x) v_x(\bigcup_{x' \neq x} E_{x'}) \sum_{x' \neq x} \mu(x') v_{x'}(\bigcup_{x' \neq x} E_{x'})]} \leq \\ &\leq \sum_{x \in X} \sqrt{[\mu(x) v_x(Y - E_x) \sum_{x' \neq x} \mu(x')]} \leq \\ &\leq \sqrt{[\text{card}(X) (1 - \min_{x' \neq x} d(v_x, v_{x'}))] \sum_{x \in X} \sqrt{[\mu(x) (1 - \mu(x))]} . \end{aligned}$$

Using this together with (2.10) and (2.11) we obtain the desired result.

**Theorem 2.** *If the loss function  $w$  is bounded from above by  $w_0$  then*

$$(2.12) \quad \frac{\gamma \mu}{2} (1 - \min_{x' \neq x} d(v_x, v_{x'})) \leq r \leq w_0 \text{card}(X) (1 - \min_{x' \neq x} d(v_x, v_{x'})) ,$$

$$(2.13) \quad \frac{\gamma \mu}{2A} (H(\mu) - I)^2 \leq r \leq \frac{w_0}{2 \log 2} (H(\mu) - I) ,$$

where  $A$  is defined by (2.7),  $\mu$  by (2.8), and  $\gamma$  by

$$\gamma = \min_{\substack{x, x' \in X \\ x' \neq x}} w(x, x') .$$

*(Left inequalities remain true also without restriction  $w \leq w_0$ .)*

**Proof.** It is clear that for every measurable disjoint decomposition  $\{E_x\}$ ,  $x \in X$ , of  $Y$  the following inequality holds

$$r \leq w_0 \sum_{x \in X} \mu(x) (1 - v_x(E_x))$$

and in order to prove the right inequality in (2.12) it remains to use Lemma 1. The left inequality immediately follows from Th. 1 in [6]. The right inequality (2.13) was proved in [6], Th. 2, and the left inequality follows from Th. 1 and from (2.12).

We shall say that a sequence  $a_n$ ,  $n = 1, 2, \dots$  of numbers converges exponentially to a number  $a$  if there exist numbers  $A > 0$  and  $0 < \lambda < 1$  such that

$$|a_n - a| < A \lambda^n \quad \text{for every } n = 1, 2, \dots$$



The loss function  $w$  will be said reflexive if  $w(x, x') > 0$  for every  $x \neq x'$ .

An immediate consequence of Theorem 2 is the following Corollary that will be very useful later.

**Corollary.** Let  $\mathscr{D}^{(n)}, v_x^{(n)}, x \in X, n = 1, 2, \dots$  be defined as in (ii) and let  $r_{(n)}$  and  $I_{(n)}$  denote the corresponding Bayes risk and information respectively. Then  $I_{(n)}$  converges to  $H(\mu)$  exponentially if and only if  $\Delta(v_x^{(n)}, v_{x'}^{(n)})$  converges to 1 exponentially for every  $x \neq x'$ . If the loss function is bounded, then this condition is sufficient in order that  $r_{(n)}$  converges to zero exponentially. If, moreover, the loss function is reflexive, then this condition is necessary and sufficient for the exponential rate of convergence of  $r_{(n)}$  to zero.

*Remark.* This corollary need not be true in case of an infinite parameter space. In order to prove this we proceed in the following manner. For ease of writing let us assume that  $X = \{1, 2, \dots\}$  and let  $w(i, j) = 0$  or 1 depending on whether  $i = j$  or  $i \neq j$ . Let  $\mathscr{D}^{(n)}$  be the  $\sigma$ -algebra of Lebesgue measurable sets in the  $n$ -dimensional Euclidian space and let  $v_i^{(n)}$  be the  $n$ -dimensional Cartesian product of uniform distribution on the interval  $\langle 2^{-i} - 1, 2^{-i} \rangle$  for every  $i \in X$ . Under this assumptions one can show that

$$\Delta(v_i^{(n)}, v_j^{(n)}) = 1 - (1 - |2^{-i} - 2^{-j}|)^n \text{ for every } n = 1, 2, \dots \text{ and } i, j \in X.$$

Hence on the hand it is clear that  $\Delta(v_i^{(n)}, v_j^{(n)})$  converges for every  $i \neq j$  exponentially to 1 as  $n \rightarrow \infty$  and on the other hand one may show on the base of Th. 1 and Th. 2 in [6] that

$$r_{(n)} \geq \frac{\mu(i)\mu(j)}{\mu(i) + \mu(j)} (1 - |2^{-i} - 2^{-j}|)^n, \quad n = 1, 2, \dots$$

$$H(\mu) - I_{(n)} \geq 2 \log 2 \frac{\mu(i)\mu(j)}{\mu(i) + \mu(j)} (1 - |2^{-i} - 2^{-j}|)^n, \quad n = 1, 2, \dots,$$

for every  $i \neq j$  so that, for every  $0 < \lambda < 1$ , there exist numbers  $\tilde{A} > 0$  and  $\tilde{\lambda} < 1$  such that

$$r_{(n)} \geq \tilde{A}\tilde{\lambda}^n,$$

$$H(\mu) - I_{(n)} \geq \tilde{A}\tilde{\lambda}^n \text{ for every } n = 1, 2, \dots$$

We are now in a position to prove by a simple contradiction that  $r_{(n)}$  does not converge to zero as well as  $I_{(n)}$  to  $H(\mu)$  exponentially for  $n \rightarrow \infty$ .

### 3. DECISION MODELS WITH INDEPENDENT AND \*-MIXING SAMPLES

In this section the classical model of statistical decision with a sample space  $(Y^n, \mathscr{D}^n)$  as it is described in Sec. 1 will be studied. We shall follow the terminology and notation employed above.

It was shown in [7] that in the independent case the condition

$$(3.1) \quad \inf_{n=1,2,\dots} \frac{1}{n} \sum_{i=1}^n A(v_{xi}, v_{x'i}) = \alpha > 0$$

implies that

$$(3.2) \quad A(v_x^n, v_{x'}^n) > 1 - 4\beta^n \quad \text{for every } n = 1, 2, \dots,$$

where  $\beta$  lies between 0 and  $e^{-\alpha/4}$ . In view of the Corollary in the preceding section it is clear that if, for every  $x \in X$ , the sequence of samples is independent and if the condition (3.1) holds for every  $x \neq x'$ , then  $I_\infty = H(\mu)$  and  $I_n$  converges to  $I_\infty$  exponentially. If moreover the loss function is bounded, then also  $r_\infty = 0$  and  $r_n$  converges to zero exponentially.

If, for every  $x \in X$ , the sequence of samples is moreover stationary, then the condition (3.1) is equivalent to the condition

$$(3.3) \quad v_x^1 \neq v_{x'}^1.$$

It follows from Th. 1 and from the considerations above, that  $I_\infty = H(\mu)$  if and only if (3.3) holds for every  $x \neq x'$ . If this condition is satisfied, then in view of (2.13) also  $r_\infty = 0$  for every bounded loss function. If, moreover, the loss function is reflexive, then in order that  $r_\infty = 0$  it is necessary and sufficient that (3.3) holds for every  $x \neq x'$ . Always when (3.3) holds  $H(\mu) - I_n$  as well as  $r_n$  converges to zero exponentially. The short discussion of the independent stationary case we can conclude by a note that the exponential convergence rate for  $I_n$  to  $H(\mu)$  under the condition that the sample space  $Y_1$  is finite was first proved by A. Rényi [5], and then generalized by author in [7]. The exponential convergence rate for  $r_n$  under the same conditions was first proved by A. Perez [3].

In Sec. 1 we have defined the model of statistical decision with  $*$ -mixing samples. According to the definition, in this case the sequence of samples is  $*$ -mixing random sequence in the sense of [1] for every realized value of the parameter. The question is which is the class of all  $*$ -mixing sequences. It is clear that independent random sequence is  $*$ -mixing. It is easy to show that if a Markov chain possesses a long-run distribution then it is  $*$ -mixing. Especially if a finite Markov chain is geometrically ergodic (in the well known sense of Kendall), then there is  $0 < \varrho < 1$  such that  $\varphi(k) = \text{const } \varrho^k$  satisfies for  $N = 1$  the requirements of the definition of  $*$ -mixing random sequence (cf. (1.5)). Hence the class of all  $*$ -mixing sequences is wider than the class of all independent sequences and consequently the  $*$ -mixing case must be studied separately.

Our considerations will be based on the following

**Lemma 2.** *If, for every  $x \in X$ , the sequence of samples is  $*$ -mixing and if, for  $i = 1, 2, \dots$ ,*

$$(3.4) \quad A(v_{xi}, v_{x'i}) \geq \alpha > 0 \quad \text{for some } x, x' \in X,$$

460 then there exists  $0 < \lambda < 1$  such that

$$(3.5) \quad A(v_x^n, v_{x'}^n) > 1 - 8\lambda^n \text{ for every } n = 1, 2, \dots$$

Proof. It follows from (3.4) that

$$\frac{1}{m} \sum_{s=1}^m A(v_{x_i s}, v_{x' i_s}) \geq \alpha$$

uniformly with respect to increasing sequences  $i_1, i_2, \dots, i_m$  of positive integers and hence, in view of the Lemma of [7], there is such  $0 < \beta < 1$  that

$$(3.6) \quad A\left(\bigotimes_{s=1}^m v_{x_i s}, \bigotimes_{s=1}^m v_{x' i_s}\right) > 1 - 4\beta^m \text{ for every } m = 1, 2, \dots$$

uniformly in the sense given above. Let  $k > N$  (cf. the definition of  $*$ -mixing random sequence) be positive integer such that

$$(3.7) \quad 1 + \varphi(k) < 1/\beta$$

and let  $n$  be an arbitrary integer. Define integers  $m$  and  $r$  by  $n = km - r$  where  $m \geq 1, 0 \leq r < k$ , and let  $\psi$  be a mapping of  $Y^n$  to  $\bigotimes_{i=1}^m Y_{ki-r}$  defined by

$$\psi((y_1, \dots, y_n)) = (y_{k-r}, y_{2k-r}, \dots, y_{mk-r}).$$

If we denote by  $\mathcal{A}^{(n)}$  the  $\sigma$ -algebra generated by the class of all sets  $E$  of the form

$$E = \bigcap_{i=1}^m \{(y_1, \dots, y_n) : y_i \in F_i\} \text{ where } F_i \in \mathcal{A}_{ki-r},$$

then the following assertions hold:

$$(a) \quad Y^{(n)} \subset Y^n,$$

(b)  $\psi$  is a measurable transformation and

$$\psi^{-1}\left(\bigotimes_{i=1}^m Y_{ki-r} - F\right) = \bigotimes_{i=1}^n Y_i - \psi^{-1}(F) \text{ for every } F \in \bigotimes_{i=1}^m \mathcal{A}_{ki-r},$$

$$(c) \quad v_x^n(\psi^{-1}(F)) \leq 1 + \varphi(k) \bigotimes_{i=1}^m v_{x_{ki-r}}(F).$$

The assertions (a) and (b) are obvious. The inequality (c) follows immediately from the definition of  $*$ -mixing sequence when  $F$  is of the form

$$F = \bigotimes_{i=1}^m F_i, \quad F_i \in \mathcal{A}_{ki-r}$$

and consequently also when  $F$  is a denumerable union of disjoint sets of this form. An application of a well-known approximation argument yields the general validity of (c).

Using (3.6) together with (i) we obtain that there exists a set

$$F \in \bigotimes_{i=1}^m \mathcal{A}_{ki-r}$$

such that

$$\begin{aligned} \bigotimes_{i=1}^m v_{xki-r} \left( \bigotimes_{i=1}^m Y_{ki-r} - F \right) &< 4\beta^m, \\ \bigotimes_{i=1}^m v_{x'ki-r} (F) &< 4\beta^m \end{aligned}$$

or, in view of (a), (b), and (c), that there exists a set  $E \in Y^n$  such that

$$\begin{aligned} v_x^n(Y^n - E) &< 4[(1 + \varphi(k)) \beta]^m, \\ v_{x'}^n(E) &< 4[(1 + \varphi(k)) \beta]^m. \end{aligned}$$

These inequalities imply in view of (i) that

$$A(v_x^n, v_{x'}^n) > 1 - 8[(1 + \varphi(k)) \beta]^m.$$

If we put  $\lambda = [(1 + \varphi(k)) \beta]^{1/k+1}$ , then by (3.7) the proof of (3.5) is complete.

**Theorem 3.** *If, for every  $x \in X$ , the sequence of samples is  $*$ -mixing and if*

$$(3.8) \quad A(v_{xi}, v_{x'i}) \geq \alpha > 0 \quad \text{for every } i = 1, 2, \dots \text{ and } x \neq x',$$

*then  $I_\infty = H(\mu)$  and  $I_n$  converges to  $I_\infty$  exponentially. If the loss function is bounded, then  $r_\infty = 0$  and  $r_n$  converges to zero exponentially.*

**Proof.** It follows from the assumptions of the Theorem that the assumptions of the preceding Lemma are satisfied for every  $x \neq x'$  and hence there exists  $0 < \lambda < 1$  such that

$$(3.9) \quad 1 - \min_{x \neq x'} A(v_x^n, v_{x'}^n) < 8\lambda^n \quad \text{for every } n = 1, 2, \dots$$

and, in accordance with (ii),  $A(v_x^\infty, v_{x'}^\infty) = 1$  for every  $x \neq x'$ . Using this together with (2.9) or (2.12) we obtain  $I_\infty = H(\mu)$  or  $r_\infty = 0$  respectively. The desired exponential rate of convergence follows from (3.9), from the Corollary in Sec. 2, and from (2.12).

We shall say that a  $*$ -mixing sequence of random samples is stationary for realized value  $x \in X$  of parameter if  $(Y_i, \mathcal{A}_i) = (Y_j, \mathcal{A}_j)$  and  $v_{xi} = v_{xj}$  for every  $i, j = 1, 2, \dots$ . It can be easily verified that this stationarity is rather weaker than the usual stationarity in the strict sense.

It is clear that if the sequence of samples is  $*$ -mixing and stationary for every  $x \in X$ , then (3.8) is satisfied if and only if  $v_x^1 \neq v_{x'}^1$  for every  $x \neq x'$  and hence we have proved the following

**Theorem 3s.** *If the sequence of samples is  $*$ -mixing and stationary for every  $x \in X$ , and if*

$$v_x^1 \neq v_{x'}^1 \text{ for every } x \neq x',$$

*then  $I_\infty = H(\mu)$  and  $I_n$  converges to  $I_\infty$  exponentially. If the loss function is bounded, then  $r_\infty = 0$  and  $r_n$  converges to zero exponentially.*

**Remark 1.** It is to be noted that the analogies between the independent and  $*$ -mixing case are not complete. Namely, a routine verification (using the Lemma of [7]) gives in the stationary independent case that the necessary and sufficient condition for the validity of  $A(v_x^\infty, v_{x'}^\infty) = 1$  is (3.3) and that if the latter condition is satisfied, then  $A(v_x^n, v_{x'}^n)$  converges to 1 exponentially as  $n \rightarrow \infty$ . An analogous assertion does not hold in the (stationary)  $*$ -mixing case. An example of statistical decision problem with stationary  $*$ -mixing samples given below shows that, even when  $v_x^1 = v_{x'}^1$ , holds, the equality  $A(v_x^\infty, v_{x'}^\infty) = 1$  as well as the exponential convergence rate of  $A(v_x^n, v_{x'}^n)$  to 1 is possible. Hence in the  $*$ -mixing case the condition (3.4) is sufficient but not necessary for the exponential rate of convergence of  $A(v_x^n, v_{x'}^n)$  to 1.

**Example 1.** Let  $X = \{x', x''\}$ ,  $Y_i = \{1, 2\}$  for every  $i = 1, 2, \dots$  and let, for every  $x$  under consideration, the sequence  $y_1, y_2, \dots$  of samples be a homogeneous Markov chain determined by an initial distribution  $v_x^1$  on  $\{1, 2\}$  and by a matrix  $W(x)$  of transition probabilities, whose element  $w_{ij}(x)$ , lying at the intersection of the  $i$ -th row and the  $j$ -th column, is given by

$$w_{ij}(x) = P[y_{n+1} = j \mid x, y_n = i] \text{ for every } n = 1, 2, \dots \text{ and } i, j = 1, 2.$$

Let us put  $v_x^1(1) = v_x^1(2) = v_{x'}^1(1) = v_{x'}^1(2) = \frac{1}{2}$ ,

$$W(x') = \begin{pmatrix} p & 1-p \\ 1-p & p \end{pmatrix},$$

$$W(x'') = \begin{pmatrix} q & 1-q \\ 1-q & q \end{pmatrix},$$

where  $p$  and  $q$  are assumed to be different. It is easily proved that in this case the sequence of samples is stationary for every  $x \in X$ . As the  $n$ -step transition probability matrices  $W^n(x')$  or  $W^n(x'')$  are given by

$$W^n(x') = \begin{pmatrix} \frac{1}{2} + (2p-1)^n & \frac{1}{2} - (2p-1)^n \\ \frac{1}{2} - (2p-1)^n & \frac{1}{2} + (2p-1)^n \end{pmatrix},$$

$$W^n(x'') = \begin{pmatrix} \frac{1}{2} + (2q-1)^n & \frac{1}{2} - (2q-1)^n \\ \frac{1}{2} - (2q-1)^n & \frac{1}{2} + (2q-1)^n \end{pmatrix}$$

respectively, it is obvious that both the homogeneous Markov chains are geometrically ergodic and consequently that the condition (1.5) is satisfied for  $N = 1$  and

$$\varphi(k) = \text{const} (\max \{|2p - 1|, |2q - 1|\})^k, \quad k = 1, 2, \dots ;$$

this proves that the sequence of samples is  $*$ -mixing for every  $x$ . On the other hand it is obvious that  $v_{x'}^1 = v_{x''}^1$ . It remains to prove that  $A(v_{x'}^n, v_{x''}^n)$  converges to 1 exponentially. Let us denote for every  $n = 1, 2, \dots$

$$Y_*^n = \bigotimes_{i=1}^n \{0, 1\}_i \quad \text{where} \quad \{0, 1\}_i = \{0, 1\}, \quad i = \overset{\circ}{\mathbb{N}}, 2, \dots$$

and let us define a mapping  $T_n$  of the space  $Y^{n+1}$  to the space  $Y_*^n$  by

$$T_n(y_1, \dots, y_{n+1}) = (y_1^*, \dots, y_n^*),$$

where

$$y_i^* = y_i \otimes y_{i+1} \quad \text{for every} \quad i = 1, 2, \dots, n,$$

and where  $y_i \otimes y_{i+1} = 1$  or 0 depending on whether  $y_i = y_{i+1}$  or  $y_i \neq y_{i+1}$ . Let us denote by  $\tilde{v}_x^n$  a probability measure induced by  $T_n$  on  $Y_*^n$ , i.e. let

$$\tilde{v}_x^n(E) = v_x^{n+1}(T_n^{-1}E) \quad \text{for every} \quad E \subset Y_*^n \quad \text{and} \quad x \in \{x', x''\}.$$

A routine verification gives that

$$\begin{aligned} \tilde{v}_x^n(y_1^*, \dots, y_n^*) &= p^{I_n}(1 - p)^{n - I_n}, \\ \tilde{v}_x^n(y_1^*, \dots, y_n^*) &= q^{I_n}(1 - q)^{n - I_n}, \end{aligned}$$

where

$$(3.10) \quad I_n = \sum_{i=1}^n y_i^*.$$

According to (iii) we get

$$A(v_{x'}^{n+1}, v_{x''}^{n+1}) \geq A(\tilde{v}_{x'}^n, \tilde{v}_{x''}^n) \geq A(B(\cdot; p, n), B(\cdot; q, n)),$$

since the distributions induced by  $\tilde{v}_{x'}^n$  or  $\tilde{v}_{x''}^n$  on the real line are  $B(\cdot; p, n)$  or  $B(\cdot; q, n)$  respectively. In view of the latter inequality, (iii), and (v), the exponential convergence rate of  $A(v_{x'}^n, v_{x''}^n)$  to 1 holds.

*Remark 2.* The question arises if the  $*$ -mixing condition is also necessary for the exponential convergence rate of  $A(v_{x'}^n, v_{x''}^n)$  to 1 (as well as of  $H(\mu) - I_n$  or  $r_n$  to zero). In the sequel we shall show that this assumption is not true. In order to achieve this we shall give the following

464 **Example 2.** Let  $X = \{x', x''\}$ ,  $Y_i = \{1, 2, 3, 4\}$ , for every  $i = 1, 2, \dots$  and let similarly as in the example above,

$$W(x') = \begin{pmatrix} p & 1-p & 0 & 0 \\ 1-p & p & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

$$W(x'') = \begin{pmatrix} q & 1-q & 0 & 0 \\ 1-q & q & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

be matrices of transition probabilities and let  $v_{x'}^1(3) = v_{x'}^1(4) = \frac{1}{2}$ ,  $v_{x'}^1(1) = v_{x'}^1(2) = v_{x''}^1(1) = v_{x''}^1(2) = \frac{1}{4}$  be initial distributions determining, for every  $x \in \{x', x''\}$ , a homogeneous Markov chain which will be referred to as a sequence of samples. It is clear that none of these random sequences is \*-mixing as it is not in this case satisfied the following necessary condition

$$\lim_{n \rightarrow \infty} \left| \frac{v_x^1(i) w_{ij}^n(x) - v_x^1(i) v_{xn}(j)}{v_x^1(i) v_{xn}(j)} \right| = 0$$

(cf. (1.5)), where  $v_{xn}(j)$ , in accordance with the notation employed above, is given by

$$v_{xn}(j) = P[y_n = j | x], \quad n = 1, 2, \dots,$$

and where  $w_{ij}^n(x)$  is an element of the  $n$ -step transition probability matrix  $W^n(x)$ , lying at the intersection of the  $i$ -th row and  $j$ -th column. In order to prove this it suffices to put, for  $x = x'$  or  $x''$ ,  $i = 2, j = 3$  or  $i = 2, j = 4$  respectively, and then to use explicit expressions for the corresponding  $n$ -step transition probability matrices.

On the other hand we shall prove that in this case  $\Delta(v_{x'}^n, v_{x''}^n)$  converges to 1 exponentially (the condition (3.8) is however in this case satisfied). It follows from the definition of the joint probability distributions  $v_x^n$ ,  $x \in X$ , and from the concrete definition of the sequence  $y_1, y_2, \dots$  of random samples above that in this special case the following equalities take place:

$$v_{x'}^n = \frac{1}{2}(v_{11}^{(n)} + v_{12}^{(n)}),$$

$$v_{x''}^n = \frac{1}{2}(v_{21}^{(n)} + v_{22}^{(n)}),$$

where  $v_{11}^{(n)}(3, 3, \dots, 3) = 1$  and  $v_{21}^{(n)}(4, 4, \dots, 4) = 1$ , and where  $v_{12}^{(n)}$  or  $v_{22}^{(n)}$  coincides with  $v_{x'}^n$  or  $v_{x''}^n$  considered in Example 1. Consequently the assumptions of (iv) in Sec. 1 are satisfied and  $\Delta(v_{12}^{(n)}, v_{22}^{(n)})$  converges to 1 exponentially. Since it is clear that  $\Delta(v_{11}^{(n)}, v_{21}^{(n)}) = 1$ , by using (iv) we conclude the proof of the desired result.

(Received December 8th, 1966.)

- [1] J. R. Blum, D. L. Hanson, L. H. Koopmans: On the strong law of large numbers for a class of stochastic processes. *Zeitschr. Wahrsch.* 2 (1963), 1.
- [2] J. L. Doob: *Stochastic processes*. J. Wiley, New York 1953.
- [3] A. Perez: Information theory methods in reducing complex decision problems. *Transactions of Fourth Prague Conf. on Inf. Theory, Stat. Dec. Functions, Random Processes*. Academia, Praha 1967.
- [4] A. Perez: Information,  $\epsilon$ -sufficiency and data reduction problems. *Kybernetika* 1 (1965), 4.
- [5] A. Rényi: On the amount of information in a sample concerning a parameter. *Publ. of Math. Inst. of Hungarian Acad. Sci.* 9 (1964), 617–625.
- [6] I. Vajda: On the statistical decision problems with discrete parameter space. *Kybernetika* 3 (1967), 2.
- [7] I. Vajda: Rate of convergence of the information in a sample concerning a parameter. *Czechoslov. Mathem. Journal* 17 (1967), 2.

## VÝTAH

---

## O statistických rozhodovacích problémech s konečným parametrovým prostorem

IGOR VAJDA

Tato práce navazuje na [6], kde byl studován klasický model statistického rozhodování s abstraktním výběrovým prostorem a s nejvýše spočetným parametrovým prostorem. Ukazuje se, že odhady základních charakteristik uvažovaného modelu, a to Bayesova rizika a střední informace o parametru obsažené ve výběru pomocí některých jednodušších veličin uvedené v [6] dávají zajímavé výsledky zejména v případě konečného parametrového prostoru.

V první části práce je definován základní model, jemu příslušné Bayesovo riziko a informace a dva speciální modely odpovídající rozhodování na základě opakovaných výběrů, které jsou v práci vyšetřeny podrobněji. V prvním z těchto modelů se předpokládá, že jednotlivé výběry jsou statisticky nezávislé, kdežto v druhém se předpokládá jistý typ slabé závislosti. V této části je dále uvedena definice  $\Delta$ -divergence dvou pravděpodobnostních distribucí a některé její vlastnosti.

V druhé části práce jsou nalezeny vztahy mezi informací resp. Bayesovým rizikem a mezi  $\Delta$ -divergencemi dvojic podmíněných pravděpodobností modelu (věty 1 a 2). Z těchto výsledků plyne, že asymptotické chování Bayesova rizika a informace (při zjemňování  $\sigma$ -algebry výběrového prostoru, speciálně při zvětšování rozsahu výběru) je dáno asymptotickým chováním příslušných  $\Delta$ -divergencí. V téže části práce je ovšem na příkladu ukázáno, že tyto závěry platí jen tehdy, když parametrový prostor je konečný.



Ve třetí části se studuje asymptotické chování  $\Delta$ -divergence podmíněných pravděpodobnostní při velkých rozsazích výběru v obou speciálních modelech definovaných výše. Ukazuje se (viz diskusi na začátku třetí části a věty 3 a 3s), že rychlost konvergence  $\Delta$ -divergence (a tedy i Bayesova rizika a informace) k jejich limitním hodnotám je v obou těchto případech za dosti obecných podmínek exponenciální.

*Ing. Igor Vajda, Ústav teorie informace a automatizace ČSAV, Praha 2, Vyšehradská 49.*