

CHARACTERIZATIONS OF PREDICTION SUFFICIENCY (ADEQUACY) IN TERMS OF RISK FUNCTIONS

BY KEI TAKEUCHI AND MASAFUMI AKAHIRA

University of Tokyo and University of Electro-Communications

Prediction sufficiency (adequacy), as it is usually defined in terms of conditional expectations, does imply "real" prediction sufficiency; i.e. sufficiency in terms of risk functions. The converse holds provided we permit the loss to depend on the unknown parameter. This is no longer true if we insist on loss functions which do not involve the unknown parameter. Conditional independence still holds but ordinary sufficiency may fail. If, however, we require equivalence of risk functions, then ordinary sufficiency and, consequently, prediction sufficiency follows.

1. Introduction. It has been shown by Bahadur [2] that sufficiency as defined in terms of conditional expectations, under regularity conditions, implies "real sufficiency" i.e. sufficiency in terms of risk functions. Furthermore it follows from Theorem 11.3 in Bahadur's paper [1] that prediction sufficiency is equivalent to ordinary sufficiency w.r.t. a larger class of probability measures. One may therefore expect similar results to hold for prediction sufficiency (adequacy) as well. In the case of prediction problems it may be of interest to consider loss functions which depend only on the decision to be made and the quantity to be predicted. If we insist on this restriction, then prediction sufficiency in terms of risk functions no longer implies prediction sufficiency as it is defined in terms of conditional expectations. Conditional independence holds but ordinary sufficiency may fail. It will, however, be shown that equivalence of risk functions implies ordinary sufficiency and consequently prediction sufficiency. (One of the authors proved this in an earlier work [7].)

We will, essentially, use the framework of Skibinsky [6]. The notion of adequacy in Skibinsky's paper is, however, replaced by the notion of prediction sufficiency.

2. Theorems. We shall assume that we are given a model consisting of a sample space $(\mathcal{X}, \mathcal{A})$ and a family $\{P_\theta : \theta \in \Theta\}$ of probability measures on \mathcal{A} . A sub σ -algebra \mathcal{B} of \mathcal{A} summarizes what can and what can not be observed. Similarly, a sub σ -algebra \mathcal{C} of \mathcal{A} describes what we are interested in predicting. Finally, we are given a sub σ -algebra \mathcal{B}_0 of \mathcal{B} and our problem is to decide if anything is lost by basing our predictions on \mathcal{B}_0 rather than \mathcal{B} .

The prediction problem is assumed to be completely described by a decision space (T, \mathcal{T}) , i.e. a measurable space and a loss function L from $\Theta \times \mathcal{X} \times T$ to $[0, \infty[$. It will always be assumed that L as a function on $\mathcal{X} \times T$ for given

Received January 1973; revised August 1974.

AMS 1970 subject classifications. Primary 62B05; Secondary 62C07.

Key words and phrases. Prediction sufficiency, conditional independence, equality of risk functions.

$\theta \in \Theta$ is $\mathcal{C} \times \mathcal{T}$ measurable. This implies that the loss does not depend on all of $x \in \mathcal{X}$, only on the part of x which is to be predicted.

A decision rule δ will here be defined as a Markov kernel $\delta(S|x): S \in \mathcal{T}$, $x \in \mathcal{X}$ which is \mathcal{B} measurable when S is fixed and a probability measure on \mathcal{T} when x is fixed.

If δ is a decision rule then its performance characteristic $\mu_\delta(\cdot|\theta)$; $\theta \in \Theta$ may be defined by defining—for each θ — $\mu_\delta(\cdot|\theta)$ as the probability measure on $\mathcal{C} \times \mathcal{T}$ defined by

$$\mu_\delta(C \times S|\theta) = \int_C \delta(S|x)P_\theta(dx).$$

The risk function $r_\delta(\theta)$; $\theta \in \Theta$ of a decision rule δ is given by:

$$r_\delta(\theta) = \int [\int L_\theta(x, t)\delta(dt|x)]P_\theta(dx), \quad \theta \in \Theta.$$

The risk function is determined by the loss function and the performance characteristic through

$$r_\delta(\theta) = \int L_\theta d\mu_\delta(\cdot|\theta), \quad \theta \in \Theta.$$

A decision rule δ will be called \mathcal{B}_0 measurable if $\delta(S|\cdot)$ is \mathcal{B}_0 measurable for each S .

DEFINITION 1. \mathcal{B} and \mathcal{C} are conditionally independent given \mathcal{B}_0 iff

(i) $P_\theta^{\mathcal{C}}(C|\mathcal{B}) = P_\theta^{\mathcal{C}}(C|\mathcal{B}_0)$ a.e. $[P_\theta]$ for all $C \in \mathcal{C}$ and for all $\theta \in \Theta$.

It is shown in Loève [5], page 351 that (i) and (ii) are equivalent:

(ii) $P_\theta^{\mathcal{A}}(B \cap C|\mathcal{B}_0) = P_\theta^{\mathcal{A}}(B|\mathcal{B}_0)P_\theta^{\mathcal{C}}(C|\mathcal{B}_0)$ a.e. $[P_\theta]$

for all $B \in \mathcal{B}$ and all $C \in \mathcal{C}$ and for all $\theta \in \Theta$.

We define prediction sufficiency (adequacy) and prediction sufficiency in the wide sense as follows:

DEFINITION 2. \mathcal{B}_0 is prediction sufficient for \mathcal{B} w.r.t. \mathcal{C} iff \mathcal{B}_0 is sufficient for \mathcal{B} and \mathcal{B} and \mathcal{C} are conditionally independent given \mathcal{B}_0 .

DEFINITION 3. \mathcal{B}_0 is prediction sufficient in the wide sense for \mathcal{B} w.r.t. \mathcal{C} iff (a) \mathcal{B} and \mathcal{C} are conditionally independent given \mathcal{B}_0 and (b) there exist \mathcal{B}_0 -measurable sets B_1 and B_2 so that $P_\theta(B_1 \cup B_2) = 1$ for all $\theta \in \Theta$ and $P_{\theta,x}^{\mathcal{A}}(\cdot|\mathcal{B}_0)$ is independent of θ if $x \in B_1$ and $P_{\theta,x}^{\mathcal{C}}(\cdot|\mathcal{B}_0)$ is independent of θ if $x \in B_2$.

In the following example we shall show that \mathcal{B}_0 is prediction sufficient in the wide sense for \mathcal{B} w.r.t. \mathcal{C} but not prediction sufficient for \mathcal{B} w.r.t. \mathcal{C} .

We assume that X_1, X_2, \dots, X_n and Y are random variables such that X_1, X_2, \dots, X_n are independently and identically distributed as $N(\theta, 1)$ while the conditional distribution of Y given X_1, X_2, \dots, X_n is $N(0, 1)$ or $N(\theta, 1)$ as $\sum_i X_i > a$ or $\sum_i X_i \leq a$. Let $\mathcal{B}, \mathcal{B}_0$ and \mathcal{C} be the σ -algebras induced by, respectively, (X_1, X_2, \dots, X_n) , $\min(a, \sum_i X_i)$ and Y . Then \mathcal{B}_0 is prediction sufficient in the wide sense for \mathcal{B} w.r.t. \mathcal{C} but not sufficient for \mathcal{B} .

That sufficiency alone is insufficient in prediction problems may be seen by considering, for example, the situation where P_θ does not depend on θ and

$\mathcal{B}_\theta = \{\phi, \mathcal{X}\}$. It is then fairly obvious that prediction of a \mathcal{C} which is not independent of \mathcal{B} should not, in general, be based on \mathcal{B}_θ .

It follows, as has been pointed out by Skibinsky [6], from Theorem 11.3 in Bahadur [1] that \mathcal{B}_θ is prediction sufficient for \mathcal{B} w.r.t. \mathcal{C} if and only if \mathcal{B}_θ is sufficient for \mathcal{B} w.r.t. all probability measures on \mathcal{B} of the form:

$$B \cap \rightarrow P_\theta(B|C),$$

where $C \in \mathcal{C}$, $\theta \in \Theta$ and $P_\theta(C) > 0$.

In analogy with Theorem 10.2 in Bahadur [1] we get:

THEOREM 1. *Suppose \mathcal{B}_θ is prediction sufficient for \mathcal{B} w.r.t. \mathcal{C} and there exists a regular conditional probability $P^\theta(\cdot | \mathcal{B}_\theta)$ of \mathcal{B} given \mathcal{B}_θ which does not depend on θ . Let δ be any decision rule from $(\mathcal{X}, \mathcal{B})$ to (T, \mathcal{T}) and put $\tilde{\delta}(S|x) = \int \delta(S|x')P_x^\theta(dx' | \mathcal{B}_\theta)$; $S \in \mathcal{T}$, $x \in \mathcal{X}$.*

Then $\tilde{\delta}$ is \mathcal{B}_θ measurable and it has the same performance characteristic as δ . In particular δ and $\tilde{\delta}$ have the same risk functions.

PROOF.

$$\begin{aligned} \mu_\theta(C \times S | \tilde{\delta}) &= \int_C \tilde{\delta}(S|x)P_\theta(dx) = \int_C E^\theta(\tilde{\delta}(S|\cdot) | \mathcal{B}_\theta) dP_\theta \\ &= \int_C E^\theta(\delta(S|\cdot) | \mathcal{B}_\theta, \mathcal{C}) dP_\theta = \int_C \delta(S|\cdot) dP_\theta = \mu_\theta(C \times S | \delta). \end{aligned}$$

REMARK. As is immediately seen from above, we need to allow randomized decision rules. This is not always necessary for the subsequent discussions.

Consequences of "risk prediction sufficiency" for various classes of loss functions. In order to show prediction sufficiency of \mathcal{B}_θ we must establish conditional independence and ordinary sufficiency. We will assume that we are given a certain class of loss functions and that to any loss function within that class and to any decision rule δ corresponds a decision rule $\tilde{\delta}$ which is \mathcal{B}_θ measurable and has uniformly smaller risk than δ . The problem is to decide whether this suffices to establish conditional independence or ordinary sufficiency. It is clear that conditional independence cannot, in general, be established by only considering loss functions which do not depend on x . Similarly loss functions which do not involve θ will, in general, be insufficient to establish ordinary sufficiency.

Conditional independence may, however, be established by considering only loss functions which do not depend on θ .

Similarly, and this follows from corresponding facts for sufficiency (see Bahadur [2], Blackwell [3] and Le Cam [4]), sufficiency of \mathcal{B}_θ may be established by considering loss functions which do not depend on x .

Conditional independence may be established by considering the two decision problem with loss functions not depending on θ as follows:

THEOREM 2. *Consider the decision space $T = \{0, 1\}$ and the set of all loss functions, L , of the form*

$$\begin{aligned} L_\theta(x, 0) &= I_C(x), & x \in \mathcal{X}, \theta \in \Theta, \\ L_\theta(x, 1) &= pI_{C^c}(x), & x \in \mathcal{X}, \theta \in \Theta, \end{aligned}$$

where $p \in]0, 1[$ and $C \in \mathcal{C}$.

Suppose that to each decision rule δ and to each loss function L , of the above form, there corresponds a \mathcal{B}_θ measurable decision rule $\tilde{\delta}$ so that

$$r_{\tilde{\delta}}(\theta) \leq r_\delta(\theta), \quad \theta \in \Theta.$$

Then \mathcal{B} and \mathcal{C} are conditionally independent given \mathcal{B}_θ .

Before proving the theorem a few remarks may be in order.

REMARK 1. It follows from the proofs that we may restrict attention to non-randomized decision rules.

REMARK 2. The proofs imply also that much smaller sets of loss functions will do. We may, for example, restrict C to a π -system generating \mathcal{C} .

REMARK 3. The parameter space Θ does not play any role in this theorem. We may—and shall—in the proof assume that Θ consists of a single point. Conditional independence is, in this situation, equivalent to prediction sufficiency.

PROOF OF THE THEOREM. We may, by Remark 3, omit the subscript θ . Furthermore a decision rule δ may be identified with the critical function $x \mapsto \delta(1|x)$. The risk may then be written:

$$\begin{aligned} r(\delta) &= \int L(\cdot, 0) dP + \int [L(\cdot, 1) - L(\cdot, 0)]\delta dP \\ &= \int L(\cdot, 0) dP + (p + 1) \int \left(\frac{P}{p + 1} - I_C \right) \delta dP \\ &= \int L(\cdot, 0) dP + (p + 1) \int \left[\frac{P}{p + 1} - P^\mathcal{C}(C|\mathcal{B}) \right] \delta dP \\ &\geq \int L(\cdot, 0) dP - (p + 1) \int \left[P^\mathcal{C}(C|\mathcal{B}) - \frac{P}{p + 1} \right]^+ dP, \end{aligned}$$

where “=” is obtained iff $\delta = 0$ a.e. or $\delta = 1$ a.e. as $P^\mathcal{C}(C|\mathcal{B}) < p/(p + 1)$ or $P^\mathcal{C}(C|\mathcal{B}) > p/(p + 1)$.

The same argument applied to \mathcal{B}_θ implies, by the assumption of the theorem, that the minimizing δ may be chosen \mathcal{B}_θ -measurable and such that $\delta = 0$ a.e. or $\delta = 1$ a.e. as $P^\mathcal{C}(C|\mathcal{B}_\theta) < p/(p + 1)$ or $P^\mathcal{C}(C|\mathcal{B}_\theta) > p/(p + 1)$. It follows that the event $[P^\mathcal{C}(C|\mathcal{B}) < p/(p + 1)]$ and the event $[P^\mathcal{C}(C|\mathcal{B}_\theta) < p/(p + 1)]$ are equivalent provided $P^\mathcal{C}(C|\mathcal{B}) \neq p/(p + 1)$ a.e. and $P^\mathcal{C}(C|\mathcal{B}_\theta) \neq p/(p + 1)$ a.e. This implies that the random variables $P^\mathcal{C}(C|\mathcal{B})$ and $P^\mathcal{C}(C|\mathcal{B}_\theta)$ have the same distribution. Hence, since $P^\mathcal{C}(C|\mathcal{B}_\theta) = E^\mathcal{C}(P(C|\mathcal{B})|\mathcal{B}_\theta)$: $P^\mathcal{C}(C|\mathcal{B}_\theta) = P^\mathcal{C}(C|\mathcal{B})$ a.e. It follows that \mathcal{B} and \mathcal{C} are conditionally independent given \mathcal{B}_θ .

REMARK. This form of the proof was suggested by one of the referees.

A criterion based on least squares prediction theory is, as has been pointed out by one of the referees, even simpler to establish. Consider a sufficiently large class of square integrable and \mathcal{C} -measurable random variables g . To a

given g we associate the loss function

$$L(x, t) = (g(x) - t)^2, \quad x \in \mathcal{X}, t \in]-\infty, \infty[.$$

Then a predictor δ minimizes the risk if and only if it is a version of $E^{\mathcal{C}}(g | \mathcal{B})$. If \mathcal{B}_θ is assumed to be just as good in this situation then $E^{\mathcal{C}}(g | \mathcal{B}) = E^{\mathcal{C}}(g | \mathcal{B}_\theta)$ a.e. This establishes conditional independence if, for example, we admit all functions $g = I_C$ where C runs through a π -system generating \mathcal{C} .

For the Lemma and Theorem 3 we assume that $\{P_\theta : \theta \in \Theta\}$ is dominated. Then we get the following lemma.

LEMMA. *If for any \mathcal{B} -measurable critical function φ there exists a \mathcal{B}_θ -measurable critical function ψ such that $E_\theta(\psi) = E_\theta(\varphi)$ for all $\theta \in \Theta$, then \mathcal{B}_θ is sufficient for \mathcal{B} .*

The proof of the lemma is essentially the same as in Bahadur [2]. The outline is as follows: Let θ_1 and θ_2 be any two points of Θ . Let ϕ be a most powerful test for θ_1 against θ_2 . Then for some k

$$\begin{aligned} \phi(x) &= 1 && \text{if } \frac{dP_{\theta_2}}{dP_{\theta_1}} > k, \\ &= 0 && \text{if } \frac{dP_{\theta_2}}{dP_{\theta_1}} < k, \end{aligned}$$

and the set $\{x : (dP_{\theta_2}/dP_{\theta_1})(x) < k\}$ is \mathcal{B}_θ -measurable. Further for every c (including ∞), the set $\{x : (dP_{\theta_2}/dP_{\theta_1})(x) < c\}$ is \mathcal{B}_θ -measurable. Hence \mathcal{B}_θ is pairwise sufficient for \mathcal{B} . Since $\{P_\theta : \theta \in \Theta\}$ is dominated, \mathcal{B}_θ is sufficient for \mathcal{B} .

The following proposition is an immediate consequence of the lemma:

Suppose that the decision space $T = [0, 1]$ and the loss function L satisfies $L(x, t) = t$, for $0 \leq t \leq 1$. If for any \mathcal{B} -measurable decision rule δ , there exists a \mathcal{B}_θ -measurable decision rule $\tilde{\delta}$ such that $r_{\tilde{\delta}}(\theta) = r_\delta(\theta)$ for all $\theta \in \Theta$, then \mathcal{B}_θ is sufficient for \mathcal{B} .

From the above, we get the following theorem.

THEOREM 3. *If for any loss function L not depending on θ and for any \mathcal{B} -measurable decision rule δ , there exists a \mathcal{B}_θ -measurable decision rule $\tilde{\delta}$ such that $r_{\tilde{\delta}}(\theta) = r_\delta(\theta)$ for all $\theta \in \Theta$, then \mathcal{B}_θ is prediction sufficient for \mathcal{B} w.r.t. \mathcal{C} .*

Various criteria for prediction sufficiency may be obtained from these results by considering, in addition to the loss functions described above, loss functions which depend on θ . If we insist on considering only loss functions which do not depend on θ then we will, in general, not be able to conclude prediction sufficiency. This follows by considering the case where $\mathcal{B}_\theta = \{\phi, \mathcal{X}\}$, \mathcal{C} is ancillary and independent of \mathcal{B} and $\{P_\theta : \theta \in \Theta\}$ is finite.

Let, in this situation, (T, \mathcal{T}) be any decision space and L any loss function which does not depend on θ and δ any decision rule from $(\mathcal{X}, \mathcal{B})$ to (T, \mathcal{T}) .

Choose a $t_\theta \in T$ such that

$$\int L(x, t_\theta) P_\theta(dx) \leq \int [\int L(x, t) P_\theta(dx)] \nu_\theta(dt), \quad \theta \in \Theta,$$

where $\nu_\theta(S) = \int \delta(S|x)P_\theta(dx)$. Then

$$r_\theta(t_\theta) \leq r_\theta(\bar{\theta}), \quad \theta \in \Theta.$$

The necessary condition for prediction sufficiency in the wide sense is obtained as follows:

THEOREM 4. *Suppose that, for each θ , there are regular conditional probabilities $P_\theta(B|\mathcal{B}_\theta): B \in \mathcal{B}$ and $P_\theta(C|\mathcal{C}_\theta): C \in \mathcal{C}$. Suppose further that \mathcal{B}_θ is prediction sufficient in the wide sense for \mathcal{B} w.r.t. \mathcal{C} .*

Let L be a loss function not depending on θ and assume that there is a \mathcal{B}_θ -measurable function τ on B_2 and an $\varepsilon \geq 0$ so that

$$\int L(x', \tau)P(dx'|\mathcal{B}_\theta) \leq \int L(x', t)P(dx'|\mathcal{B}_\theta) + \varepsilon, \quad t \in T \text{ on } B_2.$$

Then there corresponds to any decision rule δ a \mathcal{B}_θ -measurable decision rule $\bar{\delta}$ so that

$$r_\theta(\bar{\delta}) \leq r_\theta(\delta) + \varepsilon, \quad \theta \in \Theta.$$

$\bar{\delta}$ may be defined as τ on B_2 and as $E(\delta(\cdot|\cdot)|\mathcal{B}_\theta)$ on B_1 .

REMARK. There exist, for each $\varepsilon > 0$, a \mathcal{B}_θ -measurable τ on B_2 satisfying the desired inequality provided

(i) There exists a countable subset $\{t_1, t_2, \dots, t_n, \dots\}$ of T such that for all $t \in T$ and for all $x \in B_2$

$$\inf_n \int L(x', t_n)P_x^{\mathcal{C}}(dx'|\mathcal{B}_\theta) \leq \int L(x', t)P_x^{\mathcal{C}}(dx'|\mathcal{B}_\theta),$$

and

(ii) For every pair i, j and for any $\varepsilon > 0$ the set

$$M_{ij}(\varepsilon) = \{x: \int L(x', t_i)P_x^{\mathcal{C}}(dx'|\mathcal{B}_\theta) < \int L(x', t_j)P_x^{\mathcal{C}}(dx'|\mathcal{B}_\theta) + \varepsilon\}$$

is measurable.

PROOF OF THE THEOREM. Put, for each θ ,

$$\bar{\delta}_\theta^*(S|\cdot) = \int \delta(S|x')P_\theta(dx'|\mathcal{B}_\theta).$$

By the proof of Theorem 1:

$$\begin{aligned} \varepsilon + r_\theta(\bar{\delta}) &= \varepsilon + \int [\int L(x, t)\bar{\delta}_\theta^*(dt|x)]P_\theta(dx) \\ &= \int_{B_1} [\int L(x, t)\bar{\delta}(dt|x)]P_\theta(dx) \\ &\quad + \varepsilon + \int_{B_2} \{ \int [L(x', t)P_x^{\mathcal{C}}(dx'|\mathcal{B}_\theta)]\bar{\delta}_\theta^*(dt|x) \} P_\theta(dx) \\ &\geq \int_{B_1} [\int L(x, t)\bar{\delta}(dt|x)]P_\theta(dx) \\ &\quad + \int_{B_2} \{ L(x', \bar{\delta}(x))P_{\theta, x}(dx'|\mathcal{B}_\theta) \} P_\theta(dx) \\ &= \int_{B_1} [\int L(x, t)\bar{\delta}(dt|x)]P_\theta(dx) + \int_{B_2} L(x, \bar{\delta}(x))P_\theta(dx) \\ &= r_\theta(\bar{\delta}). \end{aligned}$$

Acknowledgments. The authors wish to thank Mr. M. Takahashi of Osaka University for valuable suggestions and the referees of the *Annals* for co-operation in completing the final version.

REFERENCES

- [1] BAHADUR, R. R. (1954). Sufficiency and statistical decision functions. *Ann. Math. Statist.* **25** 423-462.
- [2] BAHADUR, R. R. (1955). A characterization of sufficiency. *Ann. Math. Statist.* **26** 286-293.
- [3] BLACKWELL, D. (1953). Equivalent comparisons of experiments. *Ann. Math. Statist.* **24** 265-272.
- [4] LE CAM, L. (1964). Sufficiency and approximate sufficiency. *Ann. Math. Statist.* **35** 1419-1455.
- [5] LOÈVE, M. (1963). *Probability Theory*, (3rd. ed.). Van Nostrand, Princeton.
- [6] SKIBINSKY, M. (1967). Adequate subfields and sufficiency. *Ann. Math. Statist.* **38** 155-161.
- [7] TAKEUCHI, K. (1966). On some statistical prediction procedures. (In Japanese). *Keizaigaku Ronshu* **32** 23-31.

FACULTY OF ECONOMICS
UNIVERSITY OF TOKYO
HONGO, BUNKYO-KU
TOKYO, JAPAN

DEPARTMENT OF MATHEMATICS
UNIVERSITY OF ELECTRO-COMMUNICATIONS
CHOFUGAOKA, CHOFU-SHI
TOKYO, JAPAN