

CORRECTION

UNIVERSAL PREDICTION SCHEMES

BY PAUL H. ALGOET

The Annals of Probability (1992) **20** 901–941

In an earlier paper, the author proposed a universal scheme to learn the conditional distribution of the next outcome of a stationary process given the infinite past from past experience. An error in that paper is corrected.

1. Introduction. Let $\{X_t\}$ be a stationary random process with values in a Polish space \mathcal{X} . It follows from the martingale convergence theorem that the conditional distribution $P(dx | X^{-t})$ of $X = X_0$ given the t -past $X^{-t} = (X_{-t}, \dots, X_{-1})$ converges weakly almost surely to the conditional distribution $P(dx | X^-)$ of X given the infinite past $X^- = (\dots, X_{-2}, X_{-1})$. In Section 5 of Algoet (1992) (hereafter referred to as [A92]), we assumed that the process distribution P is unknown a priori and we constructed estimates $\hat{P}(dx | X^{-t})$ on the basis of the t -past X^{-t} such that

$$(1) \quad \hat{P}(dx | X^{-t}) \rightarrow P(dx | X^-) \text{ weakly almost surely as } t \rightarrow \infty.$$

A sequence of estimates $\hat{P}(dx | X^{-t})$ such that (1) holds under any stationary process distribution P was called a universal prediction scheme.

If \mathcal{X} is a finite set, then the prediction scheme of [A92] reduces to that of Ornstein (1978) and is universal. The scheme of [A92] was claimed to be universal in general when \mathcal{X} is a Polish space, but Gusztáv Morvai (1993; personal communication) has kindly informed me of a gap in the proof of Lemma 6. The purpose of this note is to show a way to avoid the gap. Lemma 6 may not hold as originally stated for all bounded continuous functions $h(x)$, but it does hold for certain simple functions $h(x)$, and this is all we need to prove the main result (Theorem 11) in Section 5 of [A92]. Lemma 4 was obtained by specializing Lemma 6 to the finite alphabet case, but remains valid as originally stated since Lemma 6 holds for the indicator functions $h = \delta_x$ of elements $x \in \mathcal{X}$ when \mathcal{X} is a finite set. The results on gambling, investment, modeling and data compression in the earlier sections of [A92] did not depend on Section 5 and hence are not affected.

We reformulate Lemma 6 of [A92] and introduce an approximation argument to show that the prediction scheme of [A92] is universal. We use the same notation as in [A92] with some minor changes.

Received August 1993; revised May 1994.

AMS 1991 subject classifications. 62M20, 60G25.

Key words and phrases. Learning from experience, predictive modeling, stationary processes.

2. Reformulation of Lemma 6. Let $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ be a separable metric space with its Borel σ -field and let $\{X_t\}$ be a stationary \mathcal{X} -valued random process. Let $\{\mathcal{S}_k\}_{k \geq 1}$ be an increasing sequence of finite subfields that asymptotically generate $\mathcal{B}_{\mathcal{X}}$ and let \mathcal{F}^{-t} denote the finite subfield of $\sigma(X^{-t})$ that is generated by the events $\{X_{-t} \in B_t, \dots, X_{-1} \in B_1\}$, where B_t, \dots, B_1 are \mathcal{S}_t -measurable subsets of \mathcal{X} . Notice that \mathcal{F}^{-t} is monotonically increasing and asymptotically generates the limiting σ -field $\sigma(X^-)$. For $s, t \geq 0$ we define the empirical estimate $\hat{P}_s(dx | \mathcal{F}^{-t})$ of the true conditional distribution $P(dx | \mathcal{F}^{-t})$ as in [A92].

Let $\varepsilon_k = 1/k$. We say that two distributions Q and R on the space \mathcal{X} are k -close if

$$(2) \quad \sup_{B \in \mathcal{S}_k} |Q(B) - R(B)| \leq \varepsilon_k.$$

Thus the total variation distance between the restriction of Q and the restriction of R to the finite field \mathcal{S}_k should be no larger than ε_k . If $h(x)$ is any \mathcal{S}_k -measurable function such that $0 \leq h(x) \leq 1$, then the expectations $Q\{h(X)\}$ and $R\{h(X)\}$ will differ by at most ε_k .

For each $k \geq 1$ we define the empirical estimate $\hat{P}_k(dx | X^-)$ of the true conditional distribution $P(dx | X^-)$ as follows. Find the least integer n for which there exists an integer K and a sequence $(s_i)_{0 \leq i \leq K}$ such that $k \leq K = s_0 < s_1 < \dots < s_K = n$ and all empirical estimates $\hat{P}_{s_i}(dx | \mathcal{F}^{-t})$, $s_0 \leq t \leq s_{i-1}$, $1 \leq i \leq K$ are well defined and k -close to each other. Choose K smallest possible, choose the sequence $(s_i)_{0 \leq i \leq K}$ smallest in lexicographic order when read in reverse and set

$$(3) \quad \hat{P}_k(dx | X^-) = \hat{P}_{s_K}(dx | \mathcal{F}^{-s_{K-1}}).$$

The estimate $\hat{P}_k(dx | X^-)$ was denoted by $\hat{P}(dx | X^{-\sigma_k})$ in [A92] and is well defined almost surely for all $k \geq 1$ since the search for K and $(s_i)_{0 \leq i \leq K}$ must terminate by the ergodic theorem and the martingale convergence theorem (see Lemma 5 of [A92]).

The bad event $B_{\alpha, K}^{\ell}(h)$ is defined as in [A92] for any bounded measurable function $h(x)$, any integers $K \geq 1$ and $\ell \geq 0$ and any real number $0 < \alpha < 1$. Lemma 6 of [A92] should be replaced by the following result.

LEMMA 6'. *Suppose the function $h(x)$ is \mathcal{S}_K -measurable and $0 \leq h(x) \leq 1$. Then for $\ell \geq 0$ and $0 < \alpha < 1$ we have*

$$(4) \quad P\{B_{\alpha, K}^{\ell+1}(h) | B_{\alpha, K}^{\ell}(h)\} \leq (1 - \alpha) \quad \text{almost surely}$$

and consequently, by induction on ℓ ,

$$(5) \quad P\{B_{\alpha, K}^{\ell}(h)\} \leq (1 - \alpha)^{\ell}.$$

The proof of Lemma 6 in [A92] breaks down at the foot of page 938, where it is asserted that the atom $W = W(\omega)$ is a cylinder set in $\mathcal{F}^{-\sigma_l}$ since all evidence proving that $\omega \in B_{\alpha, K}^{\ell}(h)$ is contained in $\mathcal{F}^{-\sigma_l}$. Morvai (1993) observed that

W is not $\mathcal{F}^{-\sigma_l}$ -measurable, while the assertion on page 940 that $A_E(N) \rightarrow P\{h(X) \mid \mathcal{F}^{-\sigma_l}\}$ requires that W be an atom of $\mathcal{F}^{-\sigma_l}$. [In general we have $A_E(N) \rightarrow P\{h(X) \mid W\}$.] Thus the proof of Lemma 6 is not valid for every bounded continuous or bounded measurable function $h(x)$, as was claimed. However, the lemma is valid if $h(x)$ is constant on atoms of the finite field \mathcal{S}_K , and this turns out to be sufficient for our purposes.

3. Consistency of the estimates $\hat{P}_k(dx \mid X^-)$. Theorem 11 of [A92] is valid, but the proof needs adjustment because we cannot rely on the old Lemma 6. We proceed in two steps: first we infer some conclusions from the revised Lemma 6', and then we prove universality of the estimates $\hat{P}_k(dx \mid X^-)$.

THEOREM 11A. *For any measurable set B in the generating field $\bigcup_k \mathcal{S}_k$, we have*

$$(6) \quad \lim_k \hat{P}_k\{X \in B \mid X^-\} = P\{X \in B \mid X^-\} \quad P\text{-almost surely.}$$

If a function $h(x)$ is \mathcal{S}_κ -measurable for some $\kappa \geq 1$, then the estimate $\hat{P}_k\{h(X) \mid X^-\}$ converges almost surely to the true conditional expectation $P\{h(X) \mid X^-\}$ as $k \rightarrow \infty$.

PROOF. Suppose $\kappa \geq 1$, $h(x)$ is a \mathcal{S}_κ -measurable function such that $0 \leq h(x) \leq 1$ and $0 < \alpha < 1$. We consider the event

$$(7) \quad \hat{P}_k\{h(X) \mid X^-\} \geq 2\alpha + P\{h(X) \mid X^-\} \quad \text{i.o.}$$

If this event occurs, then the event $B_{\alpha,K}^K(h)$ occurs for infinitely many K . However, if $K \geq \kappa$, then $h(x)$ is \mathcal{S}_K -measurable and Lemma 6' implies that

$$P\{B_{\alpha,K}^K(h)\} \leq (1 - \alpha)^K.$$

The event $B_{\alpha,K}^K(h)$ will occur for only finitely many K with probability 1 by the Borel–Cantelli lemma. Thus the event (7) has vanishing probability and Theorem 11A follows. \square

We now define the increasing subfields \mathcal{S}_k in a special way and prove that the estimates $\hat{P}_k(dx \mid X^-)$ converge weakly (in law, in distribution) almost surely to $P(dx \mid X^-)$ as $k \rightarrow \infty$.

Let $\{h_\kappa\}_{\kappa \geq 1}$ be a separating sequence of bounded continuous functions on \mathcal{X} . For any probability distributions Q_n and Q on \mathcal{X} , we have weak convergence $Q_n \rightarrow Q$ iff

$$Q_n\{h_\kappa(X)\} \rightarrow Q\{h_\kappa(X)\} \quad \text{as } n \rightarrow \infty \text{ for each fixed } \kappa \geq 1.$$

We assume without loss of generality that all $h_\kappa(x)$ are bounded between 0 and 1 and we define \mathcal{S}_k as the finite field that is generated by atoms of the

form

$$(8) \quad B_{i_1 i_2 \dots i_k} = \{x \in \mathcal{X} : i_\kappa 2^{-k} \leq h_\kappa(x) < (i_\kappa + 1)2^{-k}, 1 \leq \kappa \leq k\},$$

where $0 \leq i_\kappa \leq 2^k$, $1 \leq \kappa \leq k$. The functions g_1, g_2, \dots, g_k oscillate by no more than 2^{-k} on each atom $B_{i_1 i_2 \dots i_k}$ of \mathcal{S}_k . The finite fields \mathcal{S}_k are monotonically increasing and asymptotically generate the Borel σ -field $\mathcal{B}_{\mathcal{X}}$.

THEOREM 11B. *Suppose the fields \mathcal{S}_k and \mathcal{F}^{-t} are defined as above. If P is a stationary distribution, then for any bounded continuous function $h(x)$ on \mathcal{X} we have*

$$(9) \quad \lim_{k \rightarrow \infty} \hat{P}_k \{h(X) | X^-\} = P\{h(X) | X^-\} \quad P\text{-almost surely.}$$

If a regular conditional distribution $P(dx | X^-)$ exists, then

$$(10) \quad \hat{P}_k(dx | X^-) \rightarrow P(dx | X^-) \quad \text{weakly } P\text{-almost surely as } k \rightarrow \infty.$$

PROOF. It suffices to prove that for any fixed $\kappa \geq 1$ and $\varepsilon > 0$,

$$(11) \quad \limsup_k |\hat{P}_k \{h_\kappa(X) | X^-\} - P\{h_\kappa(X) | X^-\}| \leq \varepsilon \quad \text{almost surely.}$$

For each $K \geq 1$ we select a representative point $\xi_B^{(K)}$ in each atom B of \mathcal{S}_K and for any $\kappa \geq 1$ we consider the \mathcal{S}_K -measurable function

$$h_\kappa^{(K)}(x) = h_\kappa(\xi_B^{(K)}) \quad \text{if } x \in B, B \in \text{Atoms}(\mathcal{S}_K).$$

If $K \geq \kappa$, then by construction

$$(12) \quad |h_\kappa(x) - h_\kappa^{(K)}(x)| \leq 2^{-K}.$$

To prove (11), choose some K such that $K \geq \kappa$ and $2^{-K+1} \leq \varepsilon$. Obviously

$$\hat{P}_k \{h_\kappa(X) | X^-\} - P\{h_\kappa(X) | X^-\} = U_{k,\kappa}^{(K)} + V_{k,\kappa}^{(K)},$$

where

$$U_{k,\kappa}^{(K)} = \hat{P}_k \{h_\kappa^{(K)}(X) | X^-\} - P\{h_\kappa^{(K)}(X) | X^-\},$$

$$V_{k,\kappa}^{(K)} = \hat{P}_k \{h_\kappa(X) - h_\kappa^{(K)}(X) | X^-\} - P\{h_\kappa(X) - h_\kappa^{(K)}(X) | X^-\}.$$

The function $h_\kappa^{(K)}(x)$ is \mathcal{S}_K -measurable and bounded between 0 and 1; hence, $U_{k,\kappa}^{(K)} \rightarrow 0$ almost surely as $k \rightarrow \infty$ by Theorem 11A. On the other hand, it follows from (12) that $|V_{k,\kappa}^{(K)}| \leq 2^{-K+1} \leq \varepsilon$. The desired conclusion (11) and the theorem follow. \square

In general $X_t(\omega) = X(T^t \omega)$ for some random variable X with values in a separable metric space \mathcal{X} and some invertible measure-preserving transformation T of the underlying probability space (Ω, \mathcal{F}, P) . A regular conditional distribution $P(dx | X^-)$ exists if \mathcal{X} is a Polish space or a universally measurable subset of such a space, or alternatively if the measure P is perfect.

REFERENCES

- ALGOET, P. H. (1992). Universal schemes for prediction, gambling and portfolio selection. *Ann. Probab.* **20** 901–941.
- MORVAI, G. (1993). Personal communication.
- ORNSTEIN, D. S. (1978). Guessing the next output of a stationary process. *Israel J. Math.* **30** 292–296.

2613 ELMDALE
PALO ALTO, CALIFORNIA 94303