

Petr Mandl

The Kiefer-Wolfowitz approximation method in controlled Markov chains

Kybernetika, Vol. 7 (1971), No. 6, (436)--440

Persistent URL: <http://dml.cz/dmlcz/125689>

Terms of use:

© Institute of Information Theory and Automation AS CR, 1971

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>

The Kiefer-Wolfowitz Approximation Method in Controlled Markov Chains

PETR MANDL

A modification of the Kiefer-Wolfowitz stochastic approximation method is employed to maximize the mean reward per one step from a Markov chain depending on a regression parameter.

Consider a system \mathbf{S} from which income is earned at times 1, 2, 3, ... Let S_n denote the state of \mathbf{S} at time n . S_n is one of the numbers 1, 2, ..., r . The law of motion of \mathbf{S} is the following: For arbitrary $i \in \{1, 2, \dots, r\} = I$, whenever \mathbf{S} is in state i , the probability distribution of the next state is $(p_{i1}(x), \dots, p_{ir}(x))$ where $x \in (-\infty, \infty)$ is a regression parameter. The income associated with a transition from i into j equals $v_{ij}(x)$. Thus, if X_m denotes the value of the regression parameter during the period $(m, m + 1)$, then the total income earned up to time $n = 1, 2, \dots$ equals

$$V(n) = \sum_{m=1}^n v_{S_{m-1}, S_m}(X_{m-1}), \quad V(0) = 0.$$

The system is specified by matrices

$$P(x) = \|p_{ij}(x)\|_{i,j=1}^r, \quad \|v_{ij}(x)\|_{i,j=1}^r, \quad x \in (-\infty, \infty).$$

For fixed regression parameter (i.e. $X_n = x$, $n = 0, 1, \dots$), $\{S_n, n = 0, 1, \dots\}$ is a homogeneous Markov chain with transition probability matrix $P(x)$. We introduce the n -step transition probabilities $P(x)^n = \|p_{ij}^{(n)}(x)\|_{i,j=1}^r$. The expectation of $V(n)$ for $S_0 = i$ is then given by

$$E_i^x V(n) = \sum_{m=0}^{n-1} \sum_j \sum_k p_{ij}^{(m)}(x) p_{jk}(x) v_{jk}(x).$$

Assumption 1.

1. $|v_{ij}(x)| \leq K < \infty$, $x \in (-\infty, \infty)$, $i, j \in I$.

2. There exists a positive integer n_0 , an $h \in I$ and a number $d > 0$ such that

$$p_{jh}^{(n_0)}(x) \geq d, \quad j = 1, \dots, r, \quad x \in (-\infty, \infty).$$

Under Assumption 1, the limit

$$\Theta(x) = \lim_{n \rightarrow \infty} n^{-1} E_i^x V(n)$$

is independent of i . $\Theta(x)$ is the mean income per one period corresponding to regression parameter x . It can also be expressed with aid of recurrence times. Denote by $N(n)$ the n -th recurrence time into h , i.e.

$$N(0) = \inf \{m : S_m = h, m \geq 0\},$$

$$N(n) = \inf \{m : S_m = h, m > N(n-1)\}, \quad n = 1, 2, \dots$$

The pairs

$$[V(N(n+1)) - V(N(n)), N(n+1) - N(n)], \quad n = 0, 1, \dots,$$

are mutually independent, identically distributed as long as x is kept fixed. Using the strong law of large numbers it is not difficult to derive that

$$(1) \quad \Theta(x) = E_i^x [V(N(n+1)) - V(N(n))] / E_i^x [N(n+1) - N(n)].$$

We place ourselves in the situation when the dependence of Θ on x is unknown to us and we are looking for a procedure to approximate the value \hat{x} for which $\Theta(x)$ is maximal. (1) implies that we may consider this as a problem of maximizing the ratio of mean values by making independent observations on pairs of random variables. For the mean value of the ratio, i.e.

$$E_i^x \{ [V(N(n+1)) - V(N(n))] / [N(n+1) - N(n)] \},$$

the Kiefer - Wolfowitz stochastic approximation method could be applied directly. Slight modification is necessary in the present case (see Theorem 1). We shall be basing on [1] and make therefore the following assumption:

Assumption 2. $\Theta(x)$ is increasing for $x < \hat{x}$ and decreasing for $x > \hat{x}$. The derivative $\Theta'(x)$ exists and is continuous. For $x \in (-\infty, \infty)$ holds

$$K_0 |x - \hat{x}| \leq \Theta'(x) \leq K_1 |x - \hat{x}| \quad \text{where } 0 < K_0 < K_1 < \infty.$$

Description of the procedure. Let $\{a_n, n = 1, 2, \dots\}$, $\{c_n, n = 1, 2, \dots\}$ be sequences of positive numbers, $\{M_n, n = 1, 2, \dots\}$ a sequence of positive integers. Let

$$(2) \quad c_n \rightarrow 0, \quad \sum_{n=1}^{\infty} a_n = \infty, \quad \sum_{n=1}^{\infty} a_n^2 < \infty, \quad \sum_{n=1}^{\infty} a_n c_n < \infty.$$

$$(3) \quad \sum_{n=1}^{\infty} \frac{a_n}{M_n c_n} < \infty, \quad \sum_{n=1}^{\infty} \frac{a_n^2}{M_n c_n^2} < \infty.$$

Introduce $R_n = 2 \sum_{m=1}^n M_m$, $R_0 = 0$. The procedure begins by choosing an initial value x_1 of the regression parameter. At time $N(0)$ the value is altered to $x_1 + c_1$ and at time $N(M_1)$ to $x_1 - c_1$. The subsequent changes occur at times $N(R_n)$, $N(R_n + M_{n+1})$, $n = 1, 2, \dots$ in the following way: At time $N(R_n)$, x_{n+1} is calculated from

$$x_{n+1} = x_n + \left(\frac{a_n}{c_n} \right) \left[\frac{V(N(R_{n-1} + M_n)) - V(N(R_{n-1}))}{N(R_{n-1} + M_n) - N(R_{n-1})} - \frac{V(N(R_n)) - V(N(R_{n-1} + M_n))}{N(R_n) - N(R_{n-1} + M_n)} \right]$$

and the regression parameter is made equal $x_{n+1} + c_{n+1}$. At time $N(R_n + M_n)$ the parameter is altered to $x_{n+1} - c_{n+1}$. Next theorem implies that x_n converges to \hat{x} in quadratic mean.

Theorem 1. Let $\{F(y^1, y^2 | x)\}$ be a family of bivariate distribution functions depending on a real valued parameter x and such that, for an appropriate $K < \infty$,

$$\iint_{\substack{1 \leq y^1 \leq y^2 \\ |y^1| \leq K y^2}} F(dy^1, dy^2 | x) = 1, \quad \iint y^2 F(dy^1, dy^2 | x) < \infty, \quad x \in (-\infty, \infty).$$

Let the function

$$m(x) = \iint y^1 F(dy^1, dy^2 | x) / \iint y^2 F(dy^1, dy^2 | x) = m^1(x) / m^2(x)$$

be increasing for $x < \hat{x}$ and decreasing for $x > \hat{x}$. Let $m'(x)$ exist and be continuous. Assume that for each x

$$\sigma_i^2(x) = \iint (y^i - m^i(x))^2 F(dy^1, dy^2 | x) \leq \sigma^2 < \infty, \quad i = 1, 2,$$

$$K_0 |x - \hat{x}| \leq |m'(x)| \leq K_1 |x - \hat{x}|, \quad \text{where } 0 < K_0 < K_1 < \infty.$$

Let $\{a_n, n = 1, 2, \dots\}$, $\{c_n, n = 1, 2, \dots\}$ be sequences of positive numbers, $\{M_n, n = 1, 2, \dots\}$ a sequence of positive integers satisfying (2), (3). Choose x_1 arbitrary and define consecutively

$$x_{n+1} = x_n + a_n \frac{Y_{2n} - Y_{2n-1}}{c_n}, \quad n = 1, 2, \dots,$$

$$Y_{2n} = \frac{\eta_{2n,1}^1 + \eta_{2n,2}^1 + \dots + \eta_{2n,M_n}^1}{\eta_{2n,1}^2 + \eta_{2n,2}^2 + \dots + \eta_{2n,M_n}^2}, \quad Y_{2n-1} = \frac{\eta_{2n-1,1}^1 + \dots + \eta_{2n-1,M_n}^1}{\eta_{2n-1,1}^2 + \dots + \eta_{2n-1,M_n}^2},$$

and for given $\eta_{1,1}^1, \eta_{1,1}^2, \dots, \eta_{2n-2,M_{n-1}}^1, \eta_{2n-2,M_{n-1}}^2$ the vectors $(\eta_{2n-1,i}^1, \eta_{2n-1,i}^2), (\eta_{2n,i}^1, \eta_{2n,i}^2)$ $i = 1, 2, \dots, M_n$ are mutually independent with distribution function $F(y^1, y^2 | x_n - c_n)$ and $F(y^1, y^2 | x_n + c_n)$, respectively. Then

$$\lim_{n \rightarrow \infty} E(x_n - \hat{x})^2 = 0.$$

The demonstration is obtained by inserting appropriate estimates in the proof of Theorem 1 in [1] and will not be given here. Under the assumption $m''(x) \leq Q < \infty$ for $x \in (-\infty, \infty)$, it can also be shown by the methods of [1] that for

$$a_n = an^{-1}, \quad c_n = cn^{-1/4}, \quad M_n = [dn^{3/4}] + 1, \quad n = 1, 2, \dots,$$

where $a > \frac{1}{4}K_0$, $c > 0$, $d > 0$, we get

$$E(x_n - \hat{x})^2 = O(R_n^{-4/7}) \quad \text{for } n \rightarrow \infty.$$

$R_n = 2 \sum_{i=1}^n M_m$ is the number of observations employed. The corresponding estimate for the Kiefer - Wolfowitz method is

$$E(x_n - \hat{x})^2 = O(n^{-2/3}) = O(R_n^{-2/3}).$$

(Received June 3, 1971.)

REFERENCES

- [1] Václav Dupač: O Kiefer - Wolfowitzově aproximační metodě. Časopis pro pěst. mat. 82 (1957), 1, 47-75. (Appeared in Selected Translations in Mathematical Statistics and Probability.)
 [2] R. A. Howard: Dynamic Programming and Markov Processes. J. Wiley, New York 1960.

Kieferova - Wolfowitzova aproximační metoda v řízených Markovových řetězcích

PETR MANDL

V práci je modifikace Kieferovy - Wolfowitzovy stochastické aproximační metody použita k maximalizaci průměrného důchodu na jeden krok Markovova řetězce závislého na regresním parametru.

Dr. Petr Mandl, DrSc., Ústav teorie informace a automatizace ČSAV (Institute of Information Theory and Automation — Czechoslovak Academy of Sciences), Vyšehradská 49, Praha 2.