

Control of a random walk with noisy delayed information

Eitan Altman^a, Ger Koole^{b,*}

^aINRIA, Centre Sophia Antipolis, 06565 Valbonne Cedex, France

^bCWI, P.O. Box 94079, 1090 GB Amsterdam, Netherlands

Received 24 August 1993

Abstract

We consider the control of a random walk on the nonnegative integers. The controller has two actions. It makes decisions based on noisy information on the current state but on full information on previous states and actions. We establish the optimality of a threshold policy, where the threshold depends on the last action, and the noisy information. We apply the result to flow and service control problems.

Keywords: Control of a random walk; Noisy delayed information; Service control; Flow control

1. Introduction

In control problems with imperfect state information one assumes that the decision maker has only access to an observation of the state [5]. A special class of partially observable control problems is the one where information is delayed; the current state of the system becomes known to the controller only after some time T (possibly random). Such a model with T fixed was analyzed in [2]. By enlarging the state space to include the last observation of the state as well as all actions taken since that time, the model was transformed into a standard fully observable Markov decision process (MDP). Altman and Nain [2] apply the transformation to obtain an optimal control policy for a flow control model with a unit information delay. A similar flow control problem as well as a routing problem, both with a unit time of information delay, were solved in [8]. The case of N -steps delay was considered in [3]. Altman and Koole [1] analyze a problem with two controllers with delayed information on both actions and state. Koole [7] and Artiges [4] study routing with delayed information. Some other control problems with more involved information structure were considered by Hsu and Marcus [6] and Stidham [10], who study decentralized control problems. The controllers may possess some immediate noisy information about the current state of the system, yet they have access to nondelayed local information. After a unit of delay the controllers obtain the exact information. The imperfect state information problem is reduced to a perfect state information (or completely observable) problem by enlarging the state space.

In this paper we consider a control problem with a single controller with a noisy information structure, which is a special case of the information structure considered in [6]. We consider the control of a random walk on the nonnegative integers. A single controller has two actions a_0 and a_1 . It makes decisions based on

*Corresponding author.

noisy information on the current state but on full information on previous states and actions. We establish the optimality of a threshold policy, where the threshold $l(a, y)$ depends on the noisy information y on the current state and on the last action a . We then characterize the threshold, and show that $l(a_1, y) + 1 \geq l(a_0, y) \geq l(a_1, y)$. We apply the result to a problem of control of service and a problem of a control of flow in the presence of uncontrolled flows. The second application is a generalization of the model studied in [2] both in the information structure and in the more general arrival structure. The paper is organized as follows. In Section 2, we introduce the model, assumptions and notation. The main result is presented in Section 3. Extensions and applications are presented in Section 4.

2. Model and assumptions

Consider the control of a random walk in discrete time defined by $X_{n+1} = (X_n + g(\eta_n, A_n))^+$, where X_n denotes the state and A_n the action at time n . The state space X is the set of nonnegative integers, and there are two actions a_0 and a_1 available in each state. The function g is integer valued, and η_n is a sequence of \mathbb{R}^k -valued i.i.d. random variables. Further assume that the actions are ordered, i.e. $a_0 < a_1$. Note that for any $f: \mathbb{Z} \rightarrow \mathbb{R}$ such that $f(x) = f(0)$, for all $x < 0$, we have

$$E(f(X_2 - 1) | X_1 = x + 1, A_1 = a) = E(f(X_2) | X_1 = x, A_1 = a). \quad (1)$$

The applications in Section 4 consider a special case, where η_n consists of two independent components η'_n and η''_n , respectively, governing the arrivals and the departures at a queue.

In the general model, the controller does have full information on the previous states, but not on the present state. Instead, it has, at time $n + 1$, noisy information Y_n (taking values in some Borel space Y) on η_n : $Y_n = h(\eta_n, \zeta_n)$ where ζ_n is a sequence of i.i.d. random variables generating the noise. This is a special case of the information structure studied in [6]. Note that Y_n is independent of X_n and A_n . By (1), this implies that the posterior transition probabilities are also shift invariant: let $f: \mathbb{Z} \rightarrow \mathbb{R}$ be an arbitrary function such that $f(x) = f(0)$ for all $x < 0$. Then

$$E(f(X_2 - 1) | X_1 = x + 1, A_1 = a, Y_1 = y) = E(f(X_2) | X_1 = x, A_1 = a, Y_1 = y). \quad (2)$$

We model the random walk as an MDP with partial state observation. By enlarging the state space from X to $X \times A \times Y$ we obtain an equivalent fully observed MDP. Z_{n+1} , the state at time $n + 1$, is given by $Z_{n+1} = (X_n, A_n, Y_n)$.

Let Y_n have probability mass function F_2 , and let F_1 be the probability mass function of X_{n+1} , given X_n, A_n and Y_n . Thus $F_1(x|z) = P(X_n = x | Z_n = z)$ and $F_2(B) = P(Y_n \in B)$.

The transition probabilities of the MDP, for state z_n and action a_n , are

$$P(Z_{n+1} \in (\{x'\}, \{a''\}, B) | Z_n = (x, a, y), A_n = a') = 1(a'' = a') F_1(x' | (x, a, y)) F_2(B).$$

Note that (2) can be written as

$$\sum_{x' \in X} F_1(x' | x + 1, a, y) f(x' - 1) = \sum_{x' \in X} F_1(x' | x, a, y) f(x'). \quad (3)$$

We shall assume C0: $F_1(\cdot | z^2) \geq_{st} F_1(\cdot | z^1)$ for $z^i = (x^i, a^i, y^i)$, $i = 1, 2$, $y^1 = y^2$ when either $a^2 \geq a^1$ and $x^2 = x^1$, or $x^2 > x^1$. The relation \geq_{st} is the stochastic ordering, see [9]. This assumption is typical in queueing models, where the queue length is taken as the state x . C0 states intuitively that if we start with a higher (unknown) initial state we will find ourselves in a higher state after one transition. Typically, the result of having taken action a_2 instead of a_1 will result in a difference of at most one in the queue length. This is why the initial state is considered to be higher if the last known state was higher, even if the last action used is 'smaller'. Applications are presented in Section 4 that further clarify the role of assumption C0.

Consider an immediate cost $c: X \times A \rightarrow \mathbb{R}$ and assume that c satisfies C1 introduced below. For convenience, we extend the definition of c to $\mathbb{Z} \times A$ (where \mathbb{Z} is the set of integers) such that $c(y, p) = c(0, p) \forall y < 0$, $p = a_0, a_1$.

A function $f: \mathbb{Z} \times \mathcal{A} \rightarrow \mathbb{R}$ with the property that

$$f(y, p) = f(0, p) \quad \forall y < 0, p = a_0, a_1 \quad (4)$$

is said to satisfy property C1 if

- (i) $f(x, p) - f(x, q)$ is nondecreasing in x for any actions $p \geq q$.
- (ii) $f(x + 1, p) - f(x, q)$ is nondecreasing in x for any actions p and q .

It can be seen that a function $g: X \times \mathcal{A} \rightarrow \mathbb{R}$ satisfies (i), (ii) and monotonicity in x (i.e. $g(\bullet, p)$ is nondecreasing in x for any fixed p) if and only if its extension through (4) satisfies C1. Note that (ii) implies the convexity and monotone increasingness of $f(\bullet, p)$ in x for all p . This is a realistic feature of the immediate cost, which is often linear (or quadratic) in the queue size for any fixed action.

3. The optimality of a threshold policy

Let $J^n(z)$ denote the cost for a horizon of n steps, and let $J^0(z) = 0$. Define $V^n(z, a) = E(c(X_1, a) + J^n(Z_2) | Z_1 = z, A_1 = a)$. Then

$$J^{n+1}(z) = \min_a V^n(z, a). \quad (5)$$

The Markov policy u that uses in state z the action that minimizes $V^n(z, a)$ when there are n steps to go is known to be optimal. Denote $\hat{J}^n(x, a) = \int J^n(x, a, y) F_2(dy)$. We shall understand below $\hat{J}^n(x, a) = \hat{J}^n(0, a)$ for $x < 0$ (we thus consider the extension (4) of \hat{J}^n). Then $V^n(z, a') = \sum_{x'} F_1(x' | z) [c(x', a') + \hat{J}^n(x', a')]$, and (5) yields

$$\hat{J}^{n+1}(x, a) = \int F_2(dy) \left\{ \min_{a'} \sum_{x'} F_1(x' | (x, a, y)) [c(x', a') + \hat{J}^n(x', a')] \right\}.$$

Theorem 3.1. *Assume C0 and that c satisfies C1. Consider the problem of minimizing the expected cost for a horizon of n and initial state z . Then: (i) There exists an optimal policy u^* for \mathcal{Q}_n which is of a (time-dependent) threshold type such that if at time n it is optimal to use a_0 at state (x, a, y) then for any $x' > x$ it is also optimal to use a_0 at states (x', a, y) . (ii) For all $n \geq 1$, \hat{J}^n satisfies C1.*

Proof. A sufficient condition for the existence of an optimal threshold policy at stage n (when there are n steps to go) is that for $z = (x, a, y)$,

$$V^n(z, a_1) - V^n(z, a_0) \text{ is nondecreasing in } x. \quad (6)$$

Since

$$V^n(z, a_1) - V^n(z, a_0) = \sum_{x' \in X} F_1(x' | z) [c(x', a_1) - c(x', a_0) + \hat{J}^n(x', a_1) - \hat{J}^n(x', a_0)], \quad (7)$$

it follows that a sufficient condition for (6) is that both \hat{J}^n and c satisfy C1(i), and C0 holds. Indeed, this implies that the term in the square brackets of (7) is nondecreasing in x' , and (6) then follows from C0.

We show by induction that \hat{J}^n satisfies C1, hence establishing the theorem. Assume that \hat{J}^n satisfies C1. We show that for every y , $J^{n+1}(\cdot, y, \cdot)$ also satisfies C1, from which the inductive claim is established. We begin by establishing C1(i). Let p be the action that achieves the minimum in $\min_a \sum_{x' \in X} F_1(x' | x + 1, a_1, y) [c(x', a) + \hat{J}^n(x', a)]$, and let q be the action that achieves the minimum in $\min_a \sum_{x' \in X} F_1(x' | x, a_0, y) [c(x', a) + \hat{J}^n(x', a)]$. By inserting $\min_a \sum_{x' \in X} F_1(x' | \hat{x}, \hat{a}, y) [c(x', a) + \hat{J}^n(x', a)]$

for $J^{n+1}(\hat{x}, \hat{a}, y)$,

$$\begin{aligned}
& J^{n+1}(x+1, a_1, y) - J^{n+1}(x+1, a_0, y) - [J^{n+1}(x, a_1, y) - J^{n+1}(x, a_0, y)] \\
& \geq \sum_{x' \in X} F_1(x'|x+1, a_1, y)[c(x', p) + \hat{J}^n(x', p)] - \sum_{x' \in X} F_1(x'|x+1, a_0, y)[c(x', p) + \hat{J}^n(x', p)] \\
& \quad - \left[\sum_{x' \in X} F_1(x'|x+1, a_1, y)[c(x'-1, q) + \hat{J}^n(x'-1, q)] \right. \\
& \quad \quad \left. - \sum_{x' \in X} F_1(x'|x+1, a_0, y)[c(x'-1, q) + \hat{J}^n(x'-1, q)] \right] \\
& = \sum_{x' \in X} F_1(x'|x+1, a_1, y)[c(x', p) - c(x'-1, q) + \hat{J}^n(x', p) - \hat{J}^n(x'-1, q)] \\
& \quad - \sum_{x' \in X} F_1(x'|x+1, a_0, y)[c(x', p) - c(x'-1, q) + \hat{J}^n(x', p) - \hat{J}^n(x'-1, q)] \\
& \geq 0.
\end{aligned}$$

The first inequality follows from (3) and the definition of p and q . The last inequality follows from C0 and the fact that the term in square brackets is nondecreasing in x' due to C1.

It remains to establish C1(ii).

Let \hat{p} be the action that achieves the minimum in $\min_a \sum_{x' \in X} F_1(x'|x+2, p, y)[c(x', a) + \hat{J}^n(x', a)]$ and let \hat{q} be the action that achieves the minimum in $\min_a \sum_{x' \in X} F_1(x'|x, q, y)[c(x', a) + \hat{J}^n(x', a)]$. It follows from the inductive assumption that (6) holds for n and therefore $\hat{p} \leq \hat{q}$. Then, by inserting again $\min_a \sum_{x' \in X} F_1(x'|\hat{x}, \hat{a}, y)[c(x', a) + \hat{J}^n(x', a)]$ for $J^{n+1}(\hat{x}, \hat{a}, y)$,

$$\begin{aligned}
& \hat{J}^{n+1}(x+2, p, y) - \hat{J}^{n+1}(x+1, p, y) - [\hat{J}^{n+1}(x+1, q, y) - \hat{J}^{n+1}(x, q, y)] \\
& \geq \sum_{x' \in X} F_1(x'|x+2, p, y)[c(x', \hat{p}) + \hat{J}^n(x', \hat{p})] - \sum_{x' \in X} F_1(x'|x+2, p, y)[c(x'-1, \hat{q}) + \hat{J}^n(x'-1, \hat{q})] \\
& \quad - \left(\sum_{x' \in X} F_1(x'|x+1, q, y)[c(x', \hat{p}) + \hat{J}^n(x', \hat{p})] \right. \\
& \quad \quad \left. - \sum_{x' \in X} F_1(x'|x+1, q, y)[c(x'-1, \hat{q}) + \hat{J}^n(x'-1, \hat{q})] \right) \\
& = \sum_{x' \in X} F_1(x'|x+2, p, y)[c(x', \hat{p}) - c(x'-1, \hat{q}) + \hat{J}^n(x', \hat{p}) - \hat{J}^n(x'-1, \hat{q})] \\
& \quad - \sum_{x' \in X} F_1(x'|x+1, q, y)[c(x', \hat{p}) - c(x'-1, \hat{q}) + \hat{J}^n(x', \hat{p}) - \hat{J}^n(x'-1, \hat{q})] \\
& \geq 0.
\end{aligned}$$

The last inequality follows from C0, C1 and the fact that by the inductive assumption \hat{J}^n satisfies C1. This establishes the proof. \square

According to Theorem 3.1, for any a and y , there exists some threshold $l(a, y)$ such that if $x > l(a, y)$, it is optimal to use a_0 , and otherwise it is optimal to use a_1 . l satisfies the following theorem.

Theorem 3.2. For any $y \in Y$,

$$l(a_1, y) + 1 \geq l(a_0, y) \geq l(a_1, y). \quad (8)$$

Proof. Fix n . It follows from (7) that

$$\begin{aligned}
& V^n((x, a_0, y), a_0) - V^n((x, a_0, y), a_1) - [V^n((x, a_1, y), a_0) - V^n((x, a_1, y), a_1)] \\
&= \sum_{x' \in X} F_1(x' | (x, a_0, y)) [c(x', a_0) - c(x', a_1) + \hat{J}^n(x', a_0) - \hat{J}^n(x', a_1)] \\
&\quad - \sum_{x' \in X} F_1(x' | (x, a_1, y)) [c(x', a_0) - c(x', a_1) + \hat{J}^n(x', a_0) - \hat{J}^n(x', a_1)] \\
&\geq 0.
\end{aligned} \tag{9}$$

The last inequality follows from C0, and from the fact that both c and \hat{J}^n satisfy C1; hence the term in the square brackets is nonincreasing in x' .

Assume that in state $z = (x, a_0, y)$, a_0 is optimal, i.e. $V^n(z, a_1) \geq V^n(z, a_0)$. It follows from (9) that a_0 is optimal also in state (x, a_1, y) , since $V^n((x, a_1, y), a_1) \geq V^n((x, a_1, y), a_0)$. This implies the second inequality in (8). To obtain the first inequality, we note that for any $p, q \in A$,

$$\begin{aligned}
& V^n((x, p, y), a_0) - V^n((x, p, y), a_1) - [V^n((x+1, q, y), a_0) - V^n((x+1, q, y), a_1)] \\
&= \sum_{x' \in X} F_1(x' | (x, p, y)) [c(x', a_0) - c(x', a_1) + \hat{J}^n(x', a_0) - \hat{J}^n(x', a_1)] \\
&\quad - \sum_{x' \in X} F_1(x' | (x+1, q, y)) [c(x', a_0) - c(x', a_1) + \hat{J}^n(x', a_0) - \hat{J}^n(x', a_1)] \\
&\geq 0.
\end{aligned} \tag{10}$$

The last inequality follows from C0, and from the fact that both c and \hat{J}^n satisfy C1; hence the term in the square brackets is nonincreasing in x .

Assume that in state $z = (x, p, y)$, a_0 is optimal, i.e. $V^n(z, a_1) \geq V^n(z, a_0)$. It follows from (10) that a_0 is optimal also in state $(x+1, q, y)$, since $V^n((x+1, q, y), a_1) \geq V^n((x+1, q, y), a_0)$. This implies the first inequality in (8) (by setting $p = a_1$ and $q = a_0$). \square

4. Applications and extensions

The results can easily be extended to the case where ζ_n and η_n are independent but not i.i.d. Another straightforward extension is to the case where the random walk is on the set of all integers (not just the nonnegative). In that case the dynamics are given by $X_{n+1} = X_n + g(\eta_n, A_n)$. We now present applications of Theorem 3.1 to the control of queues with noisy delayed information. We then discuss situations that yield different type of noisy delayed information.

Service control

Consider a discrete-time queue with an infinite buffer. At the beginning of time n there are $\eta'_n \geq 0$ arrivals. Let X_n be the number of customers just before the arrivals occur. Let $\eta''_n \in \{0, 1\}$ be a sequence of i.i.d. Bernoulli random variables with parameter α , representing potential service completions. The action A_n corresponds to the decision whether to enable the service at time n . Let $A = \{-1, 0\}$, i.e. $a_0 = -1$ and $a_1 = 0$; take $g(\eta_n, A_n) = \eta'_n + A_n \eta''_n$. If the queue is nonempty and $A_n = a_0$ then with probability α a customer will leave the system at the end of the slot. Consider the immediate cost $c(x, a) = f(x) + \gamma a$, where f , representing a holding cost, is increasing and convex, and $\gamma \geq 0$ is a cost for deciding to serve. This structure of c ensures that C1 holds. It is easily seen that also C0 holds.

Flow control

Consider a discrete-time queue with an infinite buffer. Consider K streams of arrivals. Here $\eta_n^{(l)}$, $l = 1, \dots, K$, represents the arrival streams and $\eta_n^{(K+1)}$ represents the departures. At the beginning of time n there are $\eta_n^{(l)} \geq 0$ arrivals from sources $l = 1, \dots, K - 1$. These are uncontrolled arrivals. Let $\eta_n^{(K)} \in \{0, 1\}$ be a sequence of i.i.d. Bernoulli random variables with parameter α , representing potential arrivals from source K at the beginning of slot n . The action A_n denotes the decision of whether to enable or not the potential arrival at time n . Let $\mathcal{A} = \{0, 1\}$, i.e. $a_0 = 0$ and $a_1 = 1$; if $A_n = a_1$ then with probability α a customer will arrive from stream K . At the end of each slot, if the queue is nonempty then service succeeds with probability β and a customer leaves the system. Let $\eta_n^{(K+1)} = -1$ if service succeeds, otherwise $\eta_n^{(K+1)} = 0$. Take $g(\eta_n, A_n) = \sum_{l=1}^{K-1} \eta_n^{(l)} + A_n \eta_n^{(K)} + \eta_n^{(K+1)}$. The immediate cost is $c(x, a) = f(x) + \gamma a$, where f , representing a holding cost, is increasing and convex, and $\gamma \leq 0$ is interpreted as a reward for accepting customers, and hence a reward for increasing the throughput. This structure of c ensures that C1 holds. Again it is easily seen that C0 holds as well.

For both models we may consider the following types of information:

(1) *No information on the current state* (i.e. the state information as well as information about arrivals and service are known with a unit delay). The threshold policy obtained by Theorem 3.1 is then a function of the last action. The flow control model corresponding to this case was studied in [2] (with only one arrival stream).

(2) *Partial delayed information*. The service (or flow) controller gets in time the information from the beginning of the last slot, yet the information about events occurring at the end of the slot do not arrive in time for the decision making. Then the controller has all the information about the arrivals in the last slot but not on service completions. It thus has more information than in the previous case, but less than full information on the current state. In the case of service control this could mean that $Y_n = (\eta_n^{(l)})$, in the case of flow control $Y_n = (\eta_n^{(1)}, \dots, \eta_n^{(K)})$.

(3) *Information with a random delay*. In some cases the delay of information has random duration. In packet switching telecommunication networks, information is often obtained through acknowledgements from the destination that are piggy-backed on packets going in the opposite direction. Hence the amount of delay in the information depends on the (random) amounts of congestion of packets on the way back to the source. In our simple model, we could assume that information on the service in the last slot does not come in time for the decision making, yet with some positive probability, the information about the arrivals that occurred in the beginning of the last slot did arrive in time. In the case of the control of service this could be modeled by $Y_n = \eta_n' + \zeta_n$, where $\zeta_n = (\zeta_1, \dots, \zeta_K)$ is a vector whose components may take values 0 or $-\infty$, and are independent; $\zeta_i = 0$ means that the information comes in time, otherwise $\zeta_i = -\infty$.

(4) *Noisy delayed information*. Owing to some unreliable medium, the information we get is noisy; it becomes reliable after error correction which makes the information accurate after a unit delay.

(5) *Full information*. The controller always has full information on the current state is a special case of the information structure that we consider. Note that even for this case, our results generalize some previous results for fully observable control models. The optimality of a threshold policy for flow control models are known, e.g. [11]; the novelty of our model is that we consider the control of one flow in presence of other uncontrolled ones.

References

- [1] E. Altman and G. Koole, Stochastic scheduling games with Markov decision arrival processes, *J. Comput. Math. Appl.* **26** (6) (1993) 141–148.
- [2] E. Altman and P. Nain, Closed-loop control with delayed information, *Perf. Eval. Rev.* **20** (1992) 193–204.
- [3] E. Altman and S. Stidham, Jr., Monotonicity of optimal policies in a two-action Markov decision process, with applications to networks of queues, in preparation.
- [4] D. Artiges, Optimal routing into two heterogeneous service stations with delayed information, to appear in: *IEEE Trans. Automat. Control*.

- [5] O. Hernández-Lerma, *Adaptive Markov Control Processes* (Springer, New York, 1989).
- [6] K. Hsu and S.I. Marcus, Decentralized control of finite state Markov processes, *IEEE Trans. Automat. Control* **AC-27** (1982) 426–431
- [7] G. Koole, Optimal repairman assignment in two maintenance models which are equivalent to routing models with early decisions, Technical Report BS-R9301 CWI, Amsterdam (1993).
- [8] J. Kuri and A. Kumar, Optimal control of arrivals to queues with delayed queue length information, in: *Proc. 31st IEEE CDC, AZ, USA* (1992) 997–998.
- [9] S.M. Ross, *Stochastic Processes* (Wiley, New York, 1983).
- [10] F.C. Schoute, Decentralized control in packet switched satellite communication, *IEEE Trans. Automat. Control* **AC-23** (1978) 362–371.
- [11] S. Stidham, Jr., Optimal control of admission to a queueing system, *IEEE Trans. Automat. Control* **AC-30** (1985) 705–713.