

Tracking a Stochastic Sequence

Parisa Mansourifard, Bhaskar Krishnamachari, in Collaboration with Tara Javidi



Introduction

- We want to track a stochastic sequence, $B(t)$
- The states $i = 1, \dots, M$
- Transition matrix, P
- Partially Observable Markov Decision Process
- Finite horizon T , time steps $t = 1, \dots, T$
- Our belief vector at time t ,

$$b_t = [b_t(1), \dots, b_t(M)]$$

$$b_t(i) = \Pr(B(t) = i), i = 1, \dots, M$$

- Choose an action, a , if higher than $B(t)$, we will find out the exact amount of $B(t)$. Otherwise, we can make our belief narrower.

- Evolution of belief vector :

$$b_{t+1} = \begin{cases} T_a b_t P, & \text{if } a \leq B(t) \\ I_{B(t)} P, & \text{if } a > B(t) \end{cases}$$

$$T_a b(i) = \begin{cases} 0, & \text{if } i < a \\ \frac{b(i)}{\sum_{j=r}^M b(j)}, & \text{if } i \geq a \end{cases}$$

- Current reward :

$$R(B(t); a) = \min(a, B(t)) - C(a - B(t))^+$$

Dynamic Programming

$$V_t(b_t) = \max_{a=1, \dots, M} V_t(b_t; a), \quad \forall t = 1, \dots, T$$

$$V_T(b_T; a) = \bar{R}(b_T; a),$$

$$V_t(b_t; a) = \bar{R}(b_t; a) + \beta V_t^f(b_t; a)$$

$$V_t^f(b_t; a) = E\{V_{t+1}(b_{t+1}) | a\}$$

$$= \sum_{i=a}^M b_t(i) V_{t+1}(T_a b_t P) + \sum_{i=1}^{a-1} b_t(i) V_{t+1}(I_i P)$$

Properties of Value Function

- Convexity with respect to belief vector

$$V_t(\lambda b_1 + (1 - \lambda) b_2) \leq \lambda V_t(b_1) + (1 - \lambda) V_t(b_2)$$

$$0 \leq \lambda \leq 1$$

- Monotonically increase of future expected reward

$$V_t^f(b_t; a_1) \geq V_t^f(b_t; a_2),$$

$$a_1 \geq a_2$$

Myopic Policy

ignoring the impact of the current action on the future reward, myopic policy is given by

$$a_t^{Myopic}(b) = \arg \max_{a=1, \dots, M} \bar{R}(b_t; a)$$

$$= \min\{a = 1, \dots, M \mid \sum_{i=1}^a b_t(i) \geq \frac{1}{1+C}\}$$

Problem Formulation

- Policy vector: $\pi = [\pi(1), \dots, \pi(T)]$
- Selecting an action $\pi(t) = a_t \in \{1, \dots, M\}$
- Maximizing total discounted expected reward:

$$\max_{\pi} E^{\pi} \left[\sum_{t=1}^T \beta^{t-1} R(b_t; a_t) \mid b_1 \right]$$

- defining value function, i.e. maximum expected remaining reward starting from time t : $V_t(b)$

Optimal Policy

Our goal is to find the optimal policy or prove that it has a threshold structure.

$$a_t^{Optimal} = \arg \max_{a=1, \dots, M} V_t(b_t; a)$$

Bounds on Optimal Action:

$$a_t^{Myopic} \leq a_t^{Optimal} \leq \max\{i = 1, \dots, M \mid b_t(i) \neq 0\}$$

Future work: Find a tighter upper bound