

25 / Aug / 2019

Tabular Solution Method

Small state / action spaces.

Value function $V^*(s)$

f^n of a state

$\forall s \in S$

$Q^*(s,a)$

f^n of a state.

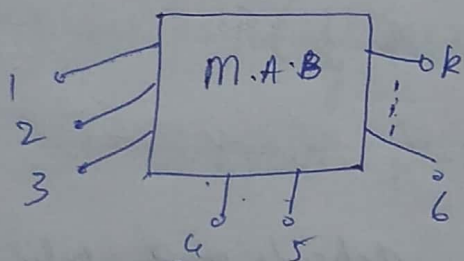
$\forall s \in S, a \in A$

Multistage problem states:

$$V^*(s) = \max_{a \in A} Q^*(s,a)$$

We estimate $V^*(s)$ & $Q^*(s,a)$. These estimates are stored in arrays.

CHAPTER - 2 MultiArm bandits.



Action: pulling the arm. —

Note: MAB has single state, but multiple action due to different rewards.

Let number of arms be 'k'

Actions available $\in \{1, 2, \dots, k\}$

Once arm is pulled, we get a reward. Given an action, i.e., follows that arm is pulled, the reward obtained follows a certain distribution. The distribution depends on the arm that is pulled.

Objective:

Find the arm that gives highest total expected reward.

Let A_t = action selected at time t

R_t = reward at time t .

Let $q_*(a) = E[R_t | A_t = a]$

Let $Q_t(a)$ = estimated value of action a at time t .
(Avg rewards)

$$Q_t(a) = \frac{\sum_{i=1}^t R_i I\{A_i = a\}}{\sum_{i=1}^t I\{A_i = a\}}$$

← total reward obtained by playing action 'a' till t .

← num of times out of t , a was played.

Where $I\{A_i = a\} = \begin{cases} 1 & \text{if } A_i = a \\ 0 & \text{otherwise} \end{cases}$

$Q_t(a)$ = Avg reward by playing 'a' till time t

We expect that

$$Q_t(a) \longrightarrow q_x(a) \text{ as time } \longrightarrow \infty$$

Example 1:

Suppose we have 5 arm ~~bandit~~ bandit m/c with actions

$$A \in \{a, b, c, d, e\}$$

Actions	R_1	R_2	R_3	R_4	R_5
a	5	6	3	2	4
b	8	4	9	2	6
c	10	12	3	1	1
d	2	3	4	8	15
e	4	3	4	2	6

Q_1	Q_2	Q_3	Q_4	Q_5
5	5.5	~5	.4	4
8	6	7	>5	<5
<u>10</u>	<u>11</u>	<u>8</u>	<u>6.5</u>	5.4
2	2.5	3	>4	<u>6.4</u>
4	3.5	<3	>3	<4

$$Q_i = \frac{\sum R_i I_{A_i=a}}{\sum I_{A_i=a}}$$

$$\text{for } a \quad Q_2(a) = \frac{R_1 I_{A_1=a} + R_2 I_{A_2=a}}{1+1} = \frac{5 \times 1 + 6 \times 1}{2} = 5.5$$

Offline Scheme: The above table represents offline

schmal. That is

5 independent copies of simulation are running (starting with same initial seed)

I simulation — arm 'a' is pulled always.

II simulation — 'b' is pulled

III simulation — 'c' is pulled.

In an online scheme, we need to decide that each time which arm to pull.

Online scheme:

Actions	R_1	R_2	R_3	R_4	R_5
a					
b		8			
c	10			12	
d			2		
e					

$$Q_t(a) = \frac{\sum_{i=1}^t R_i I_{\{A_i=a\}}}{\sum_{i=1}^t I_{\{A_i=a\}}}$$

Exploitation based scheme [Greedy approach]

$$a_t^* = \arg \max_a Q_t(a)$$

Exploration based scheme: [ϵ -greedy approach]

one eg, there are many.

$$a_t^* = \begin{cases} \arg \max_a Q_t(a) & \text{With prob } (1-\epsilon) \\ \text{Random action} & \text{With prob } \epsilon \end{cases}$$

Random action
than ~~any~~ other
than $\arg \max$

Suppose $\epsilon = 0.2$

Suppose $\arg \max_a R_t(a) = c$

ϵ -greedy scheme

Pick action c with prob 0.8

Pick action a, b, d, e with prob 0.2

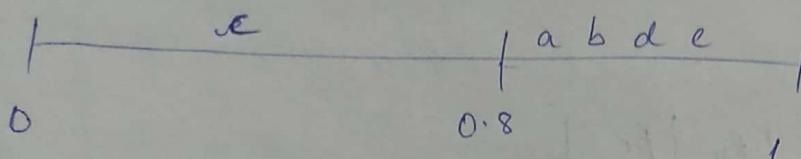


Fig 2: Ten arm, thus $k=10$. Consider 2000 randomly generated MAB. Total no of time steps in each each ~~problem~~ problem is 1000. i.e., $t=1, 2, \dots, 1000$.

Select: $q_*(a)$ according to $N(0, 1)$

$$N = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad x \in \mathbb{R}$$

↑
gaussian.

Here each sequence is played by 1000 time instances.

(Box Muller method for $N(0, 1)$ generator).

When action $A_t = a$ is picked at time t ,

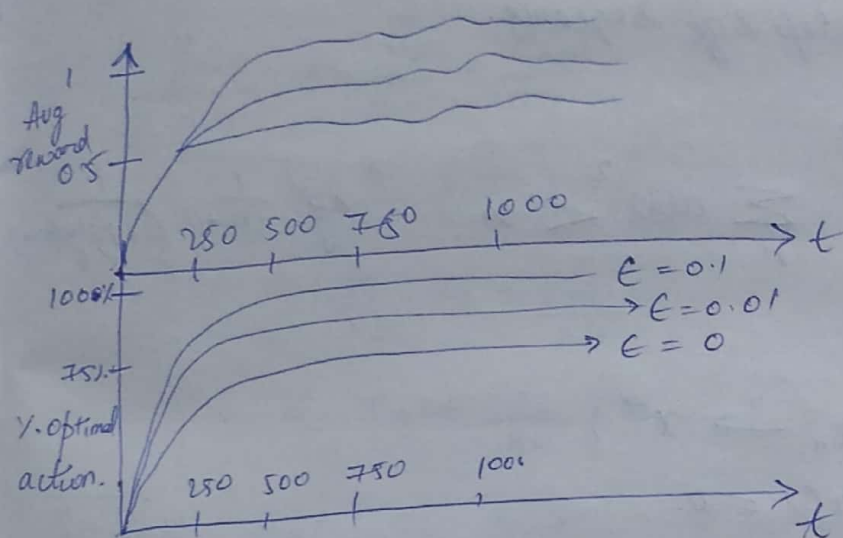
reward R_t is picked according to $N(q_*(a), 1)$

Since w.k.t, $E[Y] = E[X] + b$

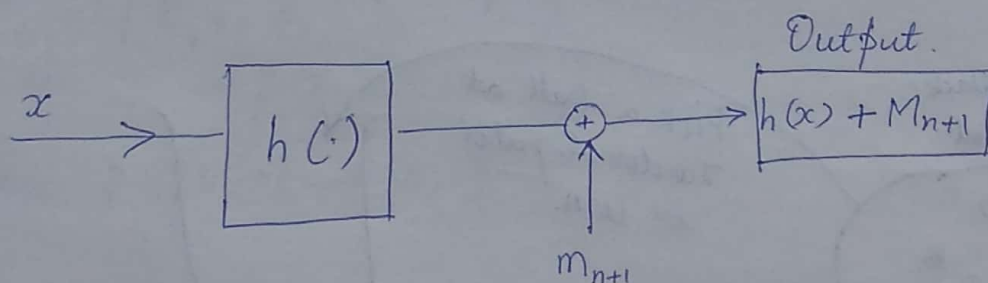
$$N(q_*(a), 1) = N(0, 1) + q_*(a)$$

$$\text{Plot } Q(a) = \frac{\sum R_i \mathbb{I}_{\{A_i=a\}}}{\sum \mathbb{I}_{\{A_i=a\}}}$$

as a function of t . (with avg over 2000 independent runs)



Stochastic Approximation Algorithms



If we input x , we observe $h(x) + M_{n+1}$ (Where M_{n+1} is zero mean noise)

Goal: Find x^* such that $h(x^*)=0$, without trying to estimate function h .

Algorithm: Start with some x_0 as an estimate of x for which $h(x) = 0$.
 (Robbins & Monro 1951)

$$x_{n+1} = x_n + a_n (h(x_n) + M_{n+1})$$

Here, $a(n)$, $n \geq 0$ is a step size sequence.

$a(n)$, $n \geq 0$ must satisfy.

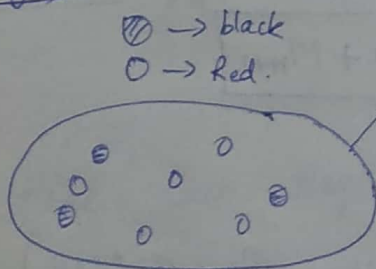
$$\sum_n a(n) = \infty \quad \sum_n a(n)^2 < \infty \quad \text{eg } \frac{1}{(n+1)} \left(\frac{1}{(n+1)^2} \right)$$

Result: As $n \rightarrow \infty$, $x_n \rightarrow x^*$,
 where $h(x^*) = 0$.

$$E[\|x_n - x^*\|^2] \rightarrow 0 \text{ as } n \rightarrow \infty \quad \left(\text{Convergence in mean sq sense} \right)$$

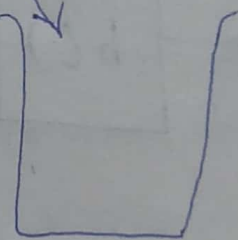
Modern day Result $x_n \rightarrow x^*$ with prob 1

Eg 3)



→ Infinite collection of B & R balls.

Pick a ball at random & put it in urn.



Initially empty Urn

Suppose: Y_n = No of Red balls in Urn by time n .

$$Y_{n+1} = Y_n + E_{n+1} \quad \text{--- ①}$$

Where $E_{n+1} = \begin{cases} 1 & \text{if red ball is picked at } n+1 \\ 0 & \text{otherwise} \end{cases}$

Define: $x_n = \frac{y_n}{n}$

ie. x_n is the fraction of red balls in urn at n .

from ①

$$\frac{y_{n+1}}{n+1} = \frac{y_n}{n+1} + \frac{\xi_{n+1}}{n+1}$$

$$= \frac{n}{n+1} \cdot \left(\frac{y_n}{n} \right) + \frac{\xi_{n+1}}{n+1}$$

$$x_{n+1} = \frac{n}{n+1} x_n + \frac{\xi_{n+1}}{n+1}$$

$$= x_n + \frac{1}{n+1} (\xi_{n+1} - x_n) \quad \text{--- (2)}$$

Suppose that a fraction of red balls in urn at time n given entire history of red & black balls that were picked. in the past depends only on x_n

$$P(\xi_{n+1} = 1 \mid \xi_0 =, \xi_1, \xi_2, \dots, \xi_n) = p(x_n)$$

Where $p: [0, 1] \rightarrow [0, 1]$

reWrite ②

$$x_{n+1} = x_n + \frac{1}{n+1} (p(x_n) - x_n) + \frac{1}{n+1} (\xi_{n+1} - p(x_n)) \quad \text{--- (3)}$$

let ~~M_{n+1}~~ $M_{n+1} \triangleq (\xi_{n+1} - p(x_n)) \Rightarrow E[M_{n+1}] = 0$

This is similar to Robbins-Monro algo

$$x_{n+1} = x_n + a_n (h(x_n) + M_{n+1})$$

Where $h(x_n) = p(x_n) - x_n$

A4 College Book

$$a_n = \frac{1}{n+1} \quad M_{n+1} = \xi_{n+1} - p(x_n)$$

Kavitha

We can consider (3) to be imitating an ODE,

$$\dot{x}(t) = f(x(t)) - x(t)$$

eg: Suppose the ODE

$$\dot{x}(t) = h(x(t)), \quad x(0) = x_0$$

$$x(t) = x(0) + \int_0^t h(x(z)) dz$$

Suppose t is small, let $t \approx \alpha$

$$x(\alpha) = x(0) + \int_0^\alpha h(x(z)) dz$$

$$\approx x(0) + \alpha h(x(\alpha))$$

assuming h doesn't vary much.

$$x((n+1)\alpha) \approx x(n\alpha) + \alpha h(x(n\alpha)) \quad n \geq 0$$

Euler's discretization.

Suppose ODE is

$$\dot{x}(t) = f(x(t)) - x(t)$$

Then Euler's discretization,

$$x((n+1)\alpha) \approx x(n\alpha) + \alpha (f(x(n\alpha)) - x(n\alpha))$$

Suppose $\hat{x}_n = x(n\alpha)$

$$\hat{x}_{n+1} = \hat{x}_n + \alpha (f(\hat{x}_n) - \hat{x}_n)$$

comparing with ③

$$x_{n+1} = x_n + \frac{1}{n+1} (p(x_n) - x_n) + \frac{1}{n+1} (\xi_{n+1} - p(x_n))$$

Suppose noise is dropped.

It turns out that,

$$\sum_n \frac{1}{n+1} (\xi_{n+1} - p(x_n)) < \infty$$

This converges!

$$Z_n = \sum_{m=0}^{n-1} \left(\frac{1}{m+1} \right) (\xi_{m+1} - p(x_m)), n \geq 1 \quad \text{is a martingale sequence.}$$

$$\text{If } E \left[\sum_n [\|Z_{n+1} - Z_n\|^2] \right] < \infty,$$

then, $\{Z_n\}$ converges.

(Quadratic Variation process)

$$\begin{aligned} (Z_{n+1} - Z_n)^2 &= \left(\frac{1}{n+1} \right)^2 (\xi_{n+1} - p(x_n))^2 \\ &= \sum_n \left(\frac{1}{n+1} \right)^2 (\xi_{n+1} - p(x_n))^2 \\ &\leq \frac{1}{n} \cdot \sum_n \frac{1}{(n+1)} < \infty \end{aligned}$$

converges.

$$x_{n+1} = x_n + \frac{1}{n+1} (p(x_n) - x_n) + \frac{1}{n+1} (\xi_{n+1} - p(x_n))$$

$$= x_0 + \sum_{m=0}^n \left(\frac{1}{m+1} \right) (p(x_m) - x_m) + \sum_{m=0}^n \left(\frac{1}{m+1} \right) (\xi_{m+1} - p(x_m))$$

$$\lim_{n \rightarrow \infty} \sum_{m=0}^n \frac{1}{m+1} (\xi_{m+1} - p(x_m)) < \infty$$

Thus, the algorithm can be considered to be a noisy Euler discretization of the ODE

$$\dot{x}(t) = p(x(t)) - x(t) \quad \text{With non uniform step size}$$

to the bounds of $\dot{x}(t)$

(I)

Suppose $x(t) = 0 \Rightarrow P(x(t)) > 0$

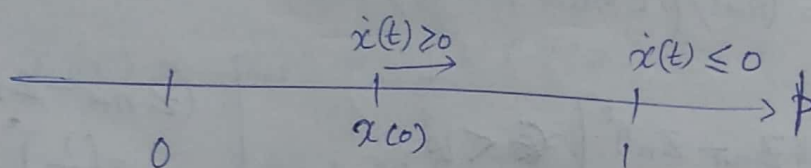
$$\therefore P(x(t)) - x(t) > 0$$

Case II

$$x(t) = 1 \Rightarrow P(x(t)) \leq 1$$

$$P(x(t)) - x(t) \leq 0$$

Suppose, ϕ is (Lipschitz Continuous)



$\dot{x}(t) > 0$ At some point in $[x(0), 1]$, $\dot{x}(t) = 0$
(Equilibrium pt of the ODE)

Suppose

$$H = \{x \mid \phi(x) = x\}$$

Then,

$$x(t) \rightarrow H \text{ as } t \rightarrow \infty$$

Then the original recursion with noise

$\{x_n\}$ satisfies

$$x_n \rightarrow H \text{ as } n \rightarrow \infty \text{ (with prob} = 1\text{)}$$

(This is not related to prev page derivation)

Estimate $Q_*(a)$

Incremental Implementation

Let $A_t = a \quad \forall t$

Let R_i = reward received after i^{th} selection action a

$$\text{Then } Q_n = \left(\frac{R_1 + R_2 + R_3 \dots + R_n}{n} \right)$$

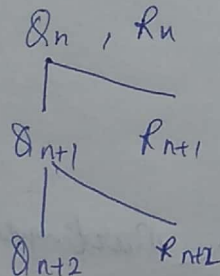
$$Q_{n+1} = \left(\frac{1}{n+1} \right) \sum_{i=1}^{n+1} R_i$$

$$= Q_n + \frac{1}{n+1} (R_{n+1} - Q_n)$$

[Incremental Implementation]

Need to store (Q_n, R_{n+1}) at any time.

this is space efficient!



Exercise: Simulate 15-arm bandit problem.

set $Q_0(a) = 5 \quad \forall a$ consider 2 cases

(a) rewards for the 15 arms \rightarrow

5, 6, 7, ..., 19

(b) rewards for 15 arms

$N(5,1), N(6,1) \dots N(19,1)$

Incorporate Greedy exploration with different values of ϵ

$$\epsilon = 0, 0.01, 0.1, 0.2$$

Plot for 5000 time steps.

II Multiple copies of this simulation 100 copies & then averages.

Consider the Increment algo again,

$$Q_{n+1} = Q_{n+1} + \frac{1}{n+1} (R_{n+1} - Q_n)$$

$$= Q_n + \frac{1}{n+1} (E[R_{n+1}] - Q_n) + \frac{1}{n+1} (R_n - E[R_{n+1}])$$

$$M_{n+1} = (R_{n+1} - E[R_{n+1}])$$

$$E[M_{n+1}] = 0.$$

Now the question is

$$\sum_{m=0}^n \left(\frac{1}{m+1} \right) (R_{m+1} - E[R_{m+1}]) < \infty ?$$

Martingale

$$\text{Suppose } |R_m| < \infty$$

$$\text{Then } \sum_{m=0}^n \left(\frac{1}{m+1} \right) (R_{m+1} - E[R_{m+1}]) < \infty$$

$$Q(t) = E[R_t] - Q(t)$$

$$\{Q^* \mid Q^* = E[R]\}$$

$$Q_n \rightarrow E[R] \text{ with prob} = 1$$

$$\frac{Q_{n+1} - Q_n}{(1/n+1)} \sim \frac{R_{n+1} - Q_n}{(E[R] - Q_n)}$$