

RL

* Exploration & Exploitation :-

5 Example :- Tic-Tac-Toe.

X	O		Player 1 & 2
	X	X	↓
		O	X O

10

Assume :-

Player 2 plays as per some strategy.

Player 1 is the RL agent.

15

* State :- $S = \{S_1, S_2, \dots, S_9, S_{10}\}$ $S_i \in \{X, O, \text{empty}\}$, $S_{10} \in \{X, O\}$

$A_i \triangleq$ change an

empty slot

by putting X or O.

20

$$S_t = (X, e, O, e, X, X, e, e, \cancel{O} \cancel{X})$$

$$S_{t+1} = (X, e, O, O, X, X, e, e, O, X)$$

25

Reward :-

$$R_i = \begin{cases} +1 & \text{if win} \\ 0 & \text{if draw} \\ -1 & \text{if lose.} \end{cases}$$

30

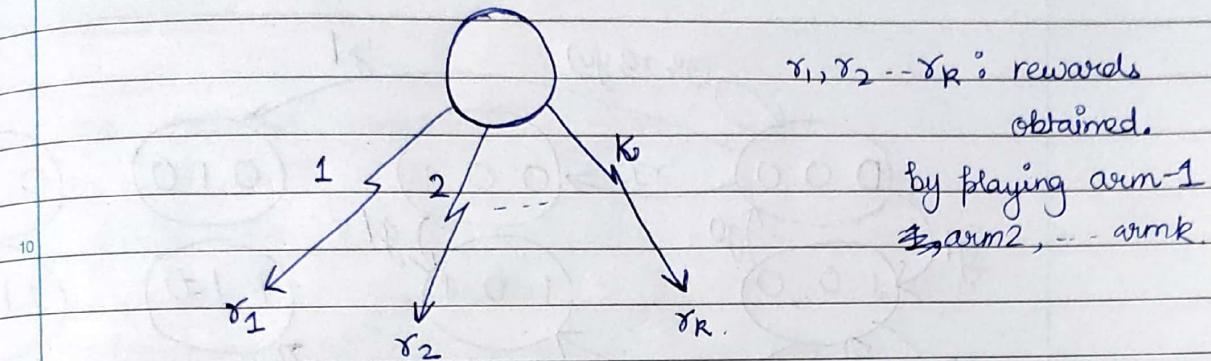
Reward :- is 0 in between (i.e. before the game ends).

Chapter-1 \rightarrow Sutton & Barto.

Chapter-2

Multi-armed Bandits

5. No. of states = 1.
multiple actions (K).



15. Assumption :- Rewards of each arm are random & come as per a certain distribution that is a function of the arm being pulled.

Problem :- find the best arm given that we don't know the reward distribution for any arm.

20. No. of states = 1.
No. of actions = K

Let

$$q^*(a) = E[R_{t+1} \mid A_t = a]$$



expected reward obtained by pulling arm a .

25. Note :- We don't know a priori what $q^*(a)$ is $\forall a=1-K$

We shall estimate $q^*(a)$ using sample average rewards

30. $Q_t(a) = \text{Sum of rewards obtained by pulling arm } a \text{ until time } t$

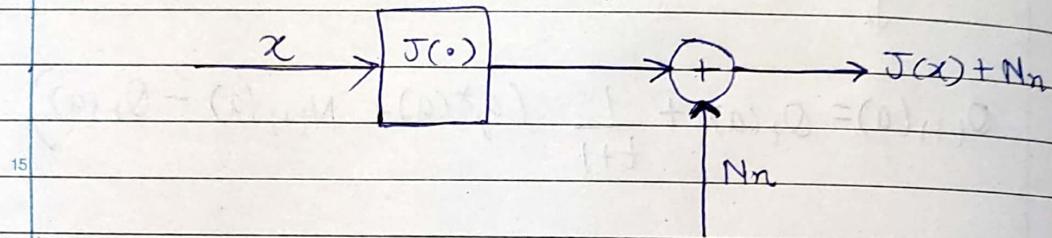
No. of time arm a is pulled until t .

New estimate = old estimate + Step-size (target -
(or learning old
rate) estimate)

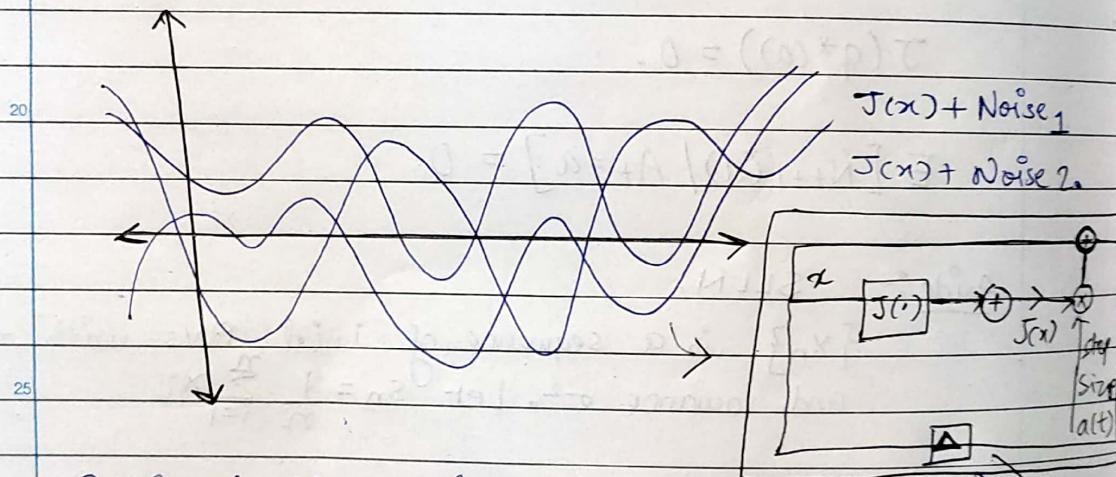
* Incremental update algos

* Stochastic Approximation algorithms:-

Suppose $J: \mathbb{R}^d \rightarrow \mathbb{R}$ is same objective function,
whose analytical form we don't know. We can obtain
noisy estimates of this function.



Q) Find x^* s.t. $J(x^*) = 0$



Robbins & Munro (1951 - Annals of Applied Stat.).

Under some assumptions, it can be shown that :-

$$x_n \rightarrow x^* \text{ w.p. 1 s.t. } J(x^*) = 0$$

Noise \rightarrow zero mean $E[N_{n+1}] = 0$

$$\begin{aligned} x_{n+1} &= x_n + \\ &\quad \alpha_n \\ &\quad [J(x_n) + N_n] \end{aligned}$$

Goal :- Find $A^* = \arg \max_a q^*(a)$.

$$Q_t(a) = \frac{\sum_{i=1}^t R_i \cdot I_{\{A_i=a\}}}{\sum_{i=1}^t I_{\{A_i=a\}}} \quad I_{\{A_i=a\}} = \begin{cases} 1 & \text{if } A_i = a \\ 0 & \text{o.w.} \end{cases}$$

Policies :-

$$(1) A_{t+1} = \arg \max_a Q_t(a)$$

(Greedy Policy) *Exploitation*

A good policy should incorporate same exploration

(2) ϵ -greedy exploration

Here, $0 < \epsilon < 1$ is a small parameter.

$$A_{t+1} = \begin{cases} \arg \max_a Q_t(a) ; \text{w.p. } (1-\epsilon). \\ \text{random action} ; \text{w.p. } \epsilon \end{cases}$$

(3) UCB (upper confidence bound) exploration :-

Suppose we pull a certain arm a all the time.

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

$$Q_{t+1}(a) = \frac{1}{(t+1)} \sum_{i=1}^{t+1} R_i = \frac{1}{(t+1)} \left(\sum_{i=1}^t R_i + R_{t+1} \right) = \frac{t}{(t+1)} \left[\frac{1}{t} \sum_{i=1}^t R_i \right] + \frac{R_{t+1}}{t+1}$$

Stochastic approx.
approx. = $\frac{t}{t+1} Q_t(a) + \frac{1}{t+1} R_{t+1}$

$$\Rightarrow Q_{t+1}(a) = Q_t(a) + \frac{1}{t+1} (R_{t+1} - Q_t(a)) \quad \text{Carmen}$$

Recall the bandit algorithm

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{t+1} (R_{t+1} - Q_t(a))$$

5

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{t+1} \left(E[R_{t+1} | A_t=a] + N_{t+1}(a) - Q_t(a) \right)$$

10

where $N_{t+1}(a) = (R_{t+1} - E[R_{t+1} | A_t=a])$



15

$$Q_{t+1}(a) = Q_t(a) + \frac{1}{t+1} (q^*(a) + N_{t+1}(a) - Q_t(a))$$

20

$$J(Q_t(a)) = (q^*(a) - Q_t(a)) \text{ is not known.}$$

$$J(q^*(a)) = 0.$$

Aside :- SLLN.

$\{X_n\}$ is a sequence of i.i.d RVs with mean μ .
and variance σ^2 . Let $S_n = \frac{1}{n} \sum_{i=1}^n X_i$.

$$S_n \rightarrow \mu \text{ a.s.}$$

30

$$S_{n+1} = S_n + \frac{1}{n+1} (X_{n+1} - S_n) \Rightarrow S_n \rightarrow \mu \text{ a.s.}$$

using Stoc. approx.

Step size or learning parameter

$$\alpha(t) = \frac{1}{t+1}, \quad t \geq 0.$$

In general, conditions on step-sizes:-

- (i) $\alpha(t) > 0 \quad \forall t.$
- (ii) $\sum_t \alpha(t) = \infty \rightarrow \text{Ensures prem}$

10

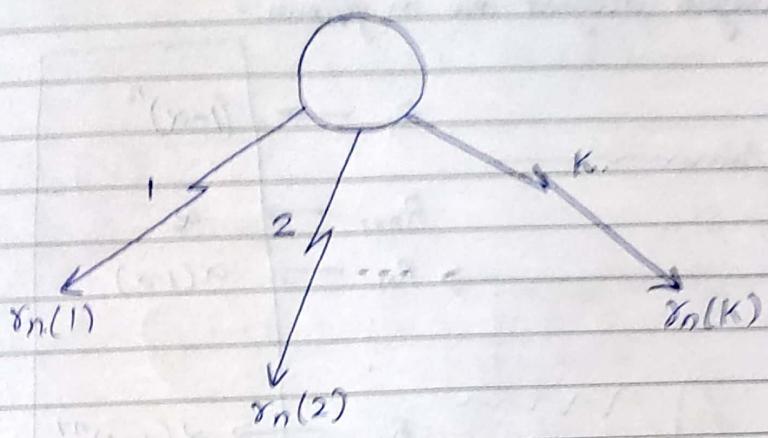
15

20

25

30

Multi-armed bandits :-



Assume $r_n(1), \dots, r_n(K)$ are independent $\forall n \neq 1, \dots, K$

So, far we assumed reward distribution does not change with time, i.e. a stationary setting.

- What happens if the dist. is changing with time, i.e. non-stationary setting.
- Const. step sizes help instead of diminishing step sizes.

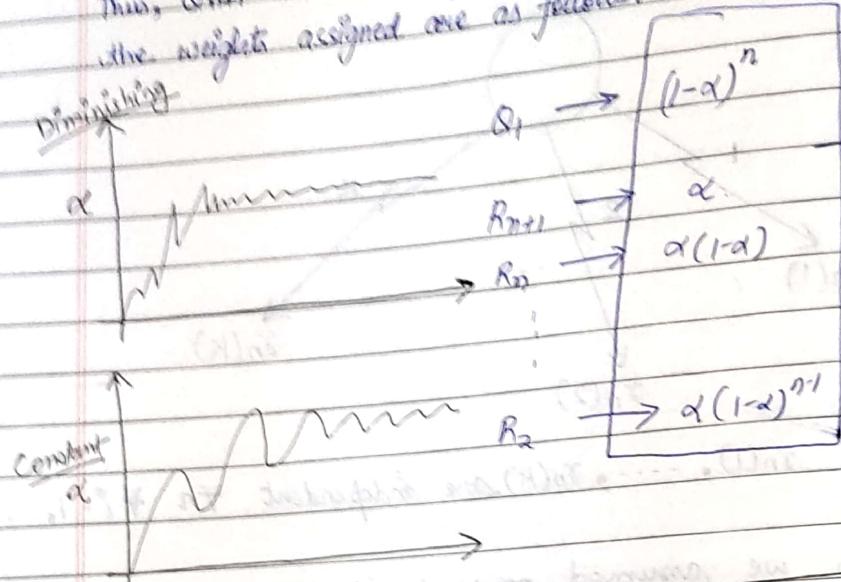
$$Q_{n+1}(i) = Q_n(i) + \alpha(R_{n+1}(i) - Q_n(i)). \text{ Here } \alpha \text{ does not change with time.}$$

Consider the recursion :-

$$\begin{aligned}
 Q_{n+1}(i) &= Q_n(i) + \alpha(R_{n+1} - Q_n) \\
 &\stackrel{\text{expand}}{=} \alpha R_{n+1} + (1-\alpha)Q_n \\
 &= \alpha R_{n+1} + (1-\alpha)(\alpha R_n + (1-\alpha)Q_{n-1}) \\
 &= \alpha R_{n+1} + \alpha(1-\alpha)R_n + (1-\alpha)^2 Q_{n-1} \\
 &\vdots \\
 &= (-\alpha)^n Q_0 + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_{i+1}
 \end{aligned}$$

Thus, $Q_{\alpha n}$ is a weighted av. of R_1, R_2, \dots, R_{n+1} , where the weights assigned are as follows:-

~~Minimizing~~



3. Exploration Strategy :- Upper Confidence Bound (UCB) :-

Let R_1, R_2, \dots, R_n are independent rewards and let them be sub-gaussian. [What is $Q_n(a)$?]

Thus, $E[R_i] = 0 \forall i$, Suppose $Q_n(a) = \frac{1}{n} \sum_{i=1}^n R_i$, $n \geq 1$.

$$P(Q_n(a) \geq \varepsilon) \leq e^{-\frac{n\varepsilon^2}{2}} = \delta.$$

$$-\frac{n\varepsilon^2}{2} = \log \delta \Rightarrow \varepsilon^2 = \frac{2 \log(\frac{1}{\delta})}{n}$$

$$\Rightarrow \varepsilon = \sqrt{\frac{2 \log(\frac{1}{\delta})}{n}}$$

Suppose $\delta = \frac{1}{n}$ (why was this arm taken?)

$$\Rightarrow P(Q_n(a) \geq \varepsilon) \leq \frac{1}{n}$$

Best plausible estimate for the maximizing arm:-

$$A_n = \arg \max_a \left(Q_n(a) + C \sqrt{\frac{\log n}{N_n(a)}} \right) \quad n \geq K$$

(Until $N_n(a) = K$, we pick each action one at a time.
 (i.e., let $N_n(a) = n$, for $n < k$).

(★) Gradient Bandit Algorithms :-

$$\text{Suppose } P(A_t = a) = \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}} \triangleq \pi_t(a).$$

(Softmax or Boltzmann distribution)

$H_t(a)$ = Preferred function.

∴ Consider the gradient algorithm :-

$$H_{t+1}(a) = H_t(a) + \alpha \frac{\partial E[R_t]}{\partial H_t(a)} ; q_*(a) = E[R_t | A_t = a]$$

$$\text{Here, } E[R_t] = \sum_x \pi_t(x) q_*(x)$$

$$\begin{aligned} \frac{\partial E[R_t]}{\partial H_t(a)} &= \frac{\partial}{\partial H_t(a)} \left[\sum_x \pi_t(x) q_*(x) \right] = \sum_x q_*(x) \frac{\partial \pi_t(x)}{\partial H_t(a)} \\ &= \sum_x (q_*(x) - B_t) \frac{\partial \pi_t(x)}{\partial H_t(a)}. \end{aligned}$$

[From where did this come].

Where B_t = Baseline function (does not depend on x).

$$\text{Note:- Extra-term :- } \sum_x B_t \frac{\partial \pi_t(x)}{\partial H_t(a)}$$

$$= B_t \sum_x \frac{\partial \pi_t(x)}{\partial H_t(a)}$$

$$= B_t \frac{\partial}{\partial H_t(a)} \left(\sum_x \pi_t(x) \right)$$

$$= 0$$

Remark: even though both estimators are unbiased, a proper choice of baseline can reduce the estimator variance.

$$\frac{\partial E[R_t]}{\partial H_t(a)} = \sum_x \pi_t(x) \left(q^*(x) - B_t \right) \frac{\frac{\partial \pi_t(x)}{\partial H_t(a)}}{\pi_t(x)}$$

$$= E \left[\left(q^*(A_t) - B_t \right) \frac{\frac{\partial \pi_t(A_t)}{\partial H_t(a)}}{\pi_t(A_t)} \right]$$

where $A_t \sim \pi_t(\cdot)$

Suppose we use $B_t = \bar{R}_t$ (av. of rewards by time t).

$$= E \left[(R_t - \bar{R}_t) \frac{\frac{\partial \pi_t(A_t)}{\partial H_t(a)}}{\pi_t(A_t)} \right]$$

$$\text{Recall } \pi_t(A_t) = \frac{e^{H_t(A_t)}}{\sum_{b=1}^K e^{H_t(b)}}$$

$$\frac{\partial \pi_t(A_t)}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left(\frac{e^{H_t(A_t)}}{\sum_{b=1}^K e^{H_t(b)}} \right)$$

$$\Rightarrow \sum_{b=1}^K c^{H_t(b)} e^{H_t(A_t)} I_{\{A_t=b\}} - e^{H_t(a)} \cdot e^{H_t(a)}$$

$$\frac{\left(\sum_{b=1}^K e^{H_t(b)} \right)^2}{}$$

$$= \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}} I_{\{A_t=a\}} \cdot \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}} - \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}}$$

(*) Algorithm

H_{t+1}

for $a = 1$

for $a \neq 1$

t

$$= \pi_t(A_t) \left(I_{\{A_t=a\}} - \pi_t(a) \right).$$

Revise \rightarrow Gradient

$$\frac{\partial E[R_t]}{\partial H_t(a)} = E \left[(R_t - \bar{R}_t) \left(I_{\{A_t=a\}} - \pi_t(a) \right) \right]$$

(★) Algorithm :-

Note: We dropped expectation. (bcz we don't know).

$$H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (I_{\{A_t=a\}} - \pi_t(a))$$

for $a = A_t$.

$$H_{t+1}(A_t) = H_t(A_t) + \alpha (R_t - \bar{R}_t) (1 - \pi_t(A_t)).$$

for $a \neq A_t$

$$H_{t+1}(a) = H_t(a) - \alpha (R_t - \bar{R}_t) \pi_t(a).$$

(★) Read ch-1 & ch-2

Sutton & Barto.

(★) Quiz on Jan 24.