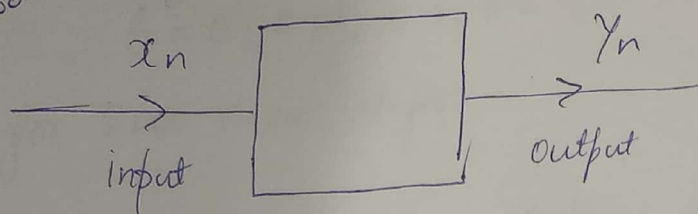


8-sept-2018)  
(1 sep notes in ~~class~~  
MBRDI-1 Book



Given  $(x_n, y_n)$ ,  $n \geq 0$  tuples & they are revealed one at a time.

Suppose  $x_n \in \mathbb{R}^n$ ,  $y_n \in \mathbb{R}^k$ ,  $n, k, \geq 0$

suppose  $f_w: \mathbb{R}^n \rightarrow \mathbb{R}^k$  where  $w = (w_1, \dots, w_d) \in \mathbb{R}^d$

$\rightarrow f_w$  is a parametrised class of functions, parameterized by  $w$ .

$\rightarrow f_w$  could be combinations of sines, coses & polynomials, NN etc

$$f_w(x) = w_0 + w_1 x + w_2 x^2 \dots + w_d x^2$$

then ~~mess~~  $w = (w_0, w_1, \dots, w_d)$

Goal: Find the best  $w$  - i.e., best explains the output as a function of input.

$$y_n = f_w(x_n) + \epsilon_n \quad \text{where } \epsilon_n \text{ is the measurement error.}$$

Now,  $\epsilon_n = y_n - f_w(x_n)$

If we are interested in min Mean Sq error, then

$$\frac{1}{2} E \|\epsilon_n\|^2 = \frac{1}{2} E \left[ \|y_n - f_w(x_n)\|^2 \right]$$

Objective  $f^u$ :

$$J(w) = \frac{1}{2} E \left[ \|y_n - f_w(x_n)\|^2 \right]$$

where  $\|\cdot\|$  is euclidean norm in  $\mathbb{R}^K$

If  $x = (x_1, x_2, x_3, \dots, x_k)$  then

$$\|x\| = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_k^2}$$

$$\|x\|^2 = x_1^2 + x_2^2 + \dots + x_n^2 \text{ i.e., } x^T x$$

$$y_n^T f_w = f_w^T y_n$$

$$\begin{aligned} \text{Then } J(w) &= \frac{1}{2} E \left[ (y_n - f_w(x_n))^T (y_n - f_w(x_n)) \right] \\ &= \frac{1}{2} E \left[ y_n^T y_n + f_w(x_n)^T f_w(x_n) - 2 f_w(x_n)^T y_n \right] \end{aligned}$$

find  $w^*$  that ~~min~~ min  $J(w)$

$$\Rightarrow \text{find } w^* = \underset{w}{\operatorname{argmin}} J(w)$$

$$\nabla_w J(w) = \frac{1}{2} E \left[ y_n^T y_n + f_n(x_n)^T f_w(x_n) - 2 f_w(x_n)^T y_n \right]$$

Assume that the gradient & expectation of gradient are inter-changeable. Then,

Note:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla f(x) = (\nabla_1 f(x), \dots, \nabla_d f(x))$$

$$\text{where } \nabla_1 f(x) = \frac{\partial f(x)}{\partial x_1}$$

$$\nabla_w J(w) = \frac{1}{2} E \left[ \nabla_w (y_n^T y_n + f_n(x_n)^T f_w(x_n) - 2 f_w(x_n)^T y_n) \right]$$

$$= \frac{1}{2} E \left[ 2 \nabla_w f_w(x_n)^T f_w(x_n) - 2 \nabla_w f_w(x_n)^T y_n \right]$$

$$= \cancel{E \nabla_w} E \left[ \nabla_w f_w(x_n) (f_w(x_n) - y_n) \right]$$

$$= - E \left[ \nabla_w f_w(x_n)^T (y_n - f_w(x_n)) \right]$$

Suppose we know the system dynamics, so the expectation is computable. Then the usual gradient descent would give

$$\begin{aligned} w_{n+1} &= w_n - \nabla_{w_n} J(w_n) \\ &= w_n + E \left[ \nabla_w f_w(x_n)^T (\gamma_n - f_{w_n}(x_n)) \right] \end{aligned}$$

Suppose we do not know the system dynamics &  $\therefore$  we don't know  $E[\cdot]$ , but we have access to samples

$$\nabla_{w_n} f_{w_n}(x_n)^T (\gamma_n - f_{w_n}(x_n))$$

Then we have,

~~$$w_{n+1} = w_n + \nabla_{w_n} f_{w_n}(x_n)^T (\gamma_n - f_{w_n}(x_n)) \quad \text{--- (1)}$$~~

$$w_{n+1} = w_n + \nabla_{w_n} f_{w_n}(x_n)^T (\gamma_n - f_{w_n}(x_n)) \quad \text{--- (1)}$$

$$w_{n+1} = w_n \quad \text{--- (2)}$$

(1) would not converge, since it iid  $\xi$  represents an random walk.

Suppose  $1 > a(n) > 0$  closer to zero

multiply (1) with  $a(n)$  & (2) with  $1 - a(n)$  and add

$$\begin{aligned} a(n)(w_{n+1}) + (1 - a(n))w_{n+1} &= a(n) \left( w_n + \nabla_{w_n} f_{w_n}(x_n)^T (\gamma_n - f_{w_n}(x_n)) \right) \\ &\quad + (1 - a(n))w_n \end{aligned}$$



~~Don't write~~

$$\Rightarrow w_{n+1} = w_n + a(n) \nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n)), n \geq 0$$

$$= w_n + a(n) E [\nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n))] + a(n) (\nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n)) - E [\nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n))])$$

$$\text{Now, } w_{n+1} = w_n - a(n) \nabla_{w_n} J(w_n) - a(n) M_{n+1}, \quad (*)$$

$$\text{where } M_{n+1} = -\nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n)) + E [\nabla_{w_n} f_{w_n}(x_n)^T (y_n - f_{w_n}(x_n))] \quad \begin{matrix} M_{n+1} \text{ is} \\ \text{mean zero} \end{matrix}$$

$$\text{We replace } \sum_n a(n) = \infty, \sum_n a(n)^2 < \infty$$

We start (t) in  $w_0$ , i.e.,

$$w(t) = w_0 + \int_0^t \nabla_w J(w(z)) dz$$

As  $t \rightarrow \infty$ ,

$$w(t) = H \equiv \{w \mid \nabla_w J(w) = 0\}$$

We can then argue that the algorithm will asymptotically converge to  $H$  with prob = 1

Stochastic Approx Scheme:

$$x_{n+1} = x_n + a(n) (h(x_n) + M_{n+1}), n \geq 0$$

~~NAME~~ ~~DATE~~ ~~TIME~~ ~~PLACE~~

~~ADITYA E RAJIV~~

~~VIKRAM NATH~~

~~KR VARUN~~

Consider a scenario where noise enters the objective as an argument  $g(x_n, \eta_n)$  where  $\eta_n$   $n \geq 0$  are i.i.d.

$$h(x_n) = E_{\eta} g(x_n, \eta_n)$$

$$M_{n+1} = g(x_n, \eta_n) - E[g(x_n, \eta_n)]$$

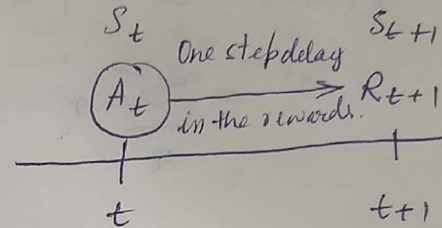
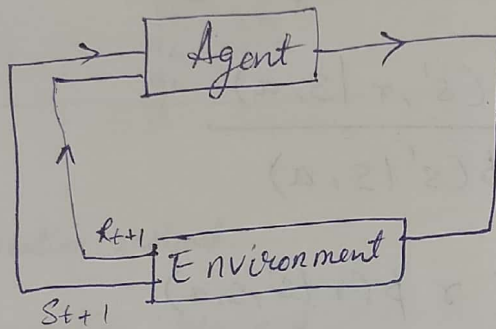
This is similar ~~to~~ ~~for~~ noise can be ~~considered as the~~ represented as previous case. thus  $g(x_n, \eta_n)$  is a special case.

# CHAPTER 3:

## Finite Markov Decision Process

Decision Maker - Agent

Agent interacts with Environment



$R_{t+1}, S_{t+1}$  depend on  $S_t, A_t$  together

The trajectory evolves as follows.

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$$

We assume that sets of states & actions & rewards are all finite.

$$\text{Let } p(s', r | s, a) \triangleq P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a) \quad \forall s, s' \in S$$

$r \in R$   
 $a \in A(s)$

Set  $S \equiv$  set of all states. In state  $s_0$ , the

set of feasible actions  $\equiv A(s)$

Now,

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1$$

$$\begin{aligned}
 p(s'|s, a) &\triangleq P(s_{t+1} = s' | s_t = s, A_t = a) \\
 &= \sum_r P(s_{t+1} = s', R_{t+1} = r | s_t = s, A_t = a) \\
 &= \sum_r p(s', r | s, a)
 \end{aligned}$$

$$\begin{aligned}
 \text{Let } r(s, a) &\triangleq E[R_{t+1} | s_t = s, A_t = a] \\
 &= \sum_r r \frac{p(s', r | s, a)}{p(s' | s, a)}
 \end{aligned}$$

$$\text{then } r(s, a) = \sum_r r p(r | s', s, a)$$

Example Movement of robotic arm to pick & place objects.  
 learning agent controls motor directly & has all recent information required. (current positions / velocity / mechanical linkages)

Actions: Voltages applied to each motor at each joint

States: Readings of joint angles & velocities.

Rewards:

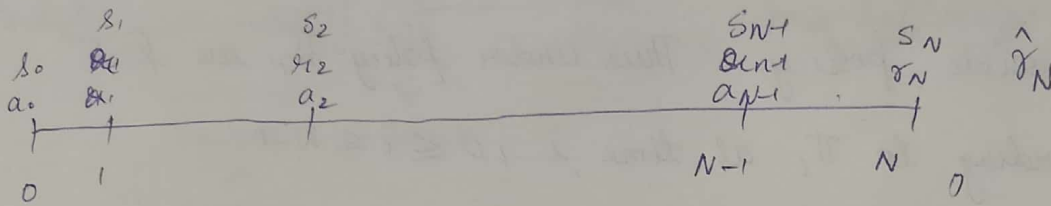
- $\begin{cases} +1 & \text{if object is correctly placed.} \\ 0 & \text{otherwise} \end{cases}$
- small -ve reward if motion is jerky.



## Exercise:

Devise three ~~examp~~ example tasks that fit into the MDP framework. Identify in each case the state, actions & rewards.

## Finite Horizon MDPs



$N$ : termination instant.

$\gamma_N$ : termination reward.

transition probabilities of states

$$p(s' | s, a)$$

$$\sum_{s'} p(s' | s, a) = 1$$

Expected single state reward

$$r_i(s, a) \triangleq E[r_{i+1} | s_i = s, A_i = a] \quad \rightarrow \text{for non terminal states.}$$

where  $i$  denotes  $i^{\text{th}}$  period.

$$\text{let } \hat{\gamma}_N(s) \triangleq E[\hat{R}_N | s_N = s] \quad \rightarrow \text{for terminal state.}$$

$\gamma_N(s)$  — expected terminal reward in state  $s$

$r_i(s, a)$  — expected single state reward



Suppose that,

$$\text{Let } \pi = \{ \pi_0, \pi_1, \dots, \pi_{N-1} \}$$

Where

$$\pi_i : S \rightarrow A \text{ such that } \pi_i(s) \in A(s) \quad \forall i=0, \dots, N-1$$

such that  $s \in S, \pi_i(s) \in A(s)$

$\pi$ : Admissible policy. Thus under policy  $\pi$ , we pick actions according to  $\pi_i$  at time  $i, 0 \leq i \leq N-1$

$$\text{Suppose } J_\pi(s_0) = E \left[ \sum_i x \right]$$

$$J_\pi(s_0) = E \left[ \sum_{i=0}^N r_i(s_i, \pi_i(s_i)) + \hat{\gamma}_N(s_N) \right]$$

Goal: Maximize  $J_\pi(s_0)$  wrt  $\pi$

Find  $\pi^*$  s.t.

$$J^*(s_0) \triangleq \max_{\pi} J_\pi(s_0) = J_{\pi^*}(s_0)$$

Result:  $J^*(s_0)$  can be obtained by following algorithm.

$$V_N(s_N) = \hat{\gamma}_N(s_N)$$

$$V_k(s_k) = \max_{a_k} \sum_k p(s_k, a_k, s_{k+1}) \left( r_k(s_k, a_k) + V_{k+1}(s_{k+1}) \right) \hat{\gamma}_N(s_N)$$

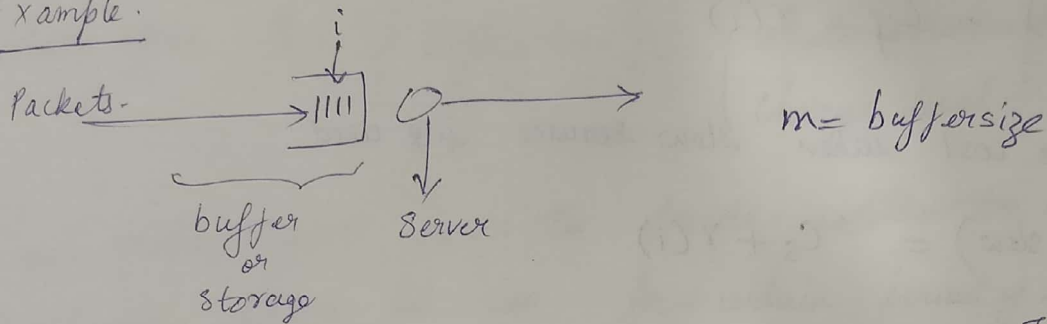
$$k = N, N-1, \dots, 1, 0$$

Joining back words in time.

$$V_0(s_0) = J^*(s_0)$$

in the above result,  $s_0$  is not zero state, it is applicable to all the states.

Example:



Server can provide fast service or slow service. Fast service comes at a cost  $C_f$ . Slow service comes at a cost  $C_s$  such that  $C_f > C_s$ .

let  $P_m$  = prob of  $m$  arrivals in a slot.

- A packet can take 1 or more slots for service. and ends service just before end of a slot.
- If fast service is used, prob of service completion by the end of slot =  $P_f$
- If slow service is used, prob of service completion by the end of slot =  $P_s$

$$P_s > P_f$$

Darwin

sunny  
+

sy sunny  
+

sunny  
+

Suppose  $h(i)$  = cost of holding 'i' packets in a system

Where i is no of packets in the system.

- Single stage cost when ~~the~~ fast server is used

$$g(i, \text{fast}) = C_f + r(i)$$

- Single stage cost when slow service is used

$$g(i, \text{slow}) = C_s + r(i)$$

$$p(m|0, \text{fast}) = \sum_{n=m}^{\infty} p_n$$

0 ~~at the~~ given that ~~of~~ zero packets initially.

$$p(m+1) - p(m|1, \text{fast}) = \text{Case 1: Customer in service leaves system by end of slot.}$$

$p_f$

$$(1-p_f) \sum_{i=m-1}^{\infty} p_i$$

Case 2: Customer does not leave

(m-1) arrivals  $(1-p_f)$