

Finite horizon Problems :-

By a controlled Markov chain we mean a sequence of random variables on a common probability space $\{X_n\}$ taking values in the same set S and whose dynamics depends on a control sequence $\{Z_n\}$ st. $Z_n \in A \quad \forall n$ for some set A .

* Controlled Markov property :-

$$P(X_{n+1} = j | X_n = i, Z_n = z, X_{n-1} = i_{n-1}, Z_{n-1} = z_{n-1}) = P(X_{n+1} = j | X_n = i, Z_n = z) \triangleq p(i, z, j)$$

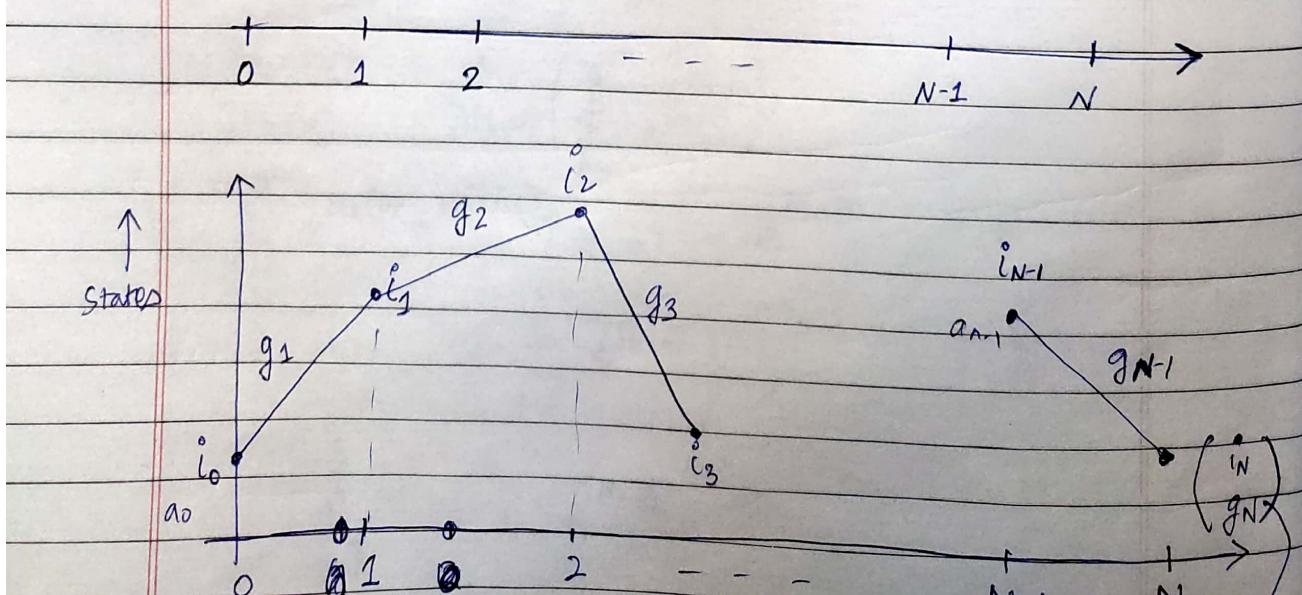
* Markov Decision Process :- (MDP) :-

A Markov Decision Process (MDP) is a controlled Markov process with an additional cost structure.

- A cost $g(i, z, j)$ is incurred when

$X_n = i, Z_n = z, X_{n-1} = j$

Let N be the horizon, $N < \infty$



* Long term cost objective :-

$$J_{\{z_0, z_1, \dots, z_{n-1}\}}(i_0) = E \left[g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), x_{k+1}) \right] \quad | \quad x_0 = i_0$$

Here, $E[\cdot]$ is the expectation over the joint distribution of the next states, x_1, x_2, \dots, x_N given that $x_0 = i_0$.

* Policy :- Decision rule specified is a set of functions that suggests the action to chose at a given instant.

* Let $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$ where $\mu_k : S \rightarrow A$ $\forall k = 0, 1, \dots, N-1$

* We say that π is an admissible policy if

$$\mu_k : S \rightarrow A$$

$\mu_k(i) \rightarrow A(i)$ [set of feasible actions in state i]

$$A(i) \subset A \quad \forall i \in S \quad \bigcup_{i \in S} A(i) = A$$

* Deterministic

policy.

* Stochastic

Policies.

* Optimal Policy :- An optimal policy π^* is one for which

$$\cancel{J_{\pi^*}(i_0)} = \min_{\pi \in \Pi} J_{\pi}(i_0) \quad \forall i_0 \in S.$$

where $\Pi =$ set of all admissible policies.

Here,

$$J_{\pi}(i_0) = E \left[g(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), x_{k+1}) \right] \quad | \quad x_0 = i_0$$

$$\text{Here } \pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$$

We also refer to the optimal cost as

$$J^*(i_0) = J_{\pi^*}(i_0)$$

optimal
cost

Optimal policy.

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{argmin}} J_{\pi}(i_0)$$

Let $\pi^* = \{\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*\}$ be the optimal policy. Assume that under π^* , a state x_i occurs at time i with the positive probability.
Consider the following subproblem:

$$\pi^i = \{m_i, m_{i+1}, \dots, m_{N-1}\}$$

$$E[g_N(x_N) + \sum_{k=i}^{N-1} g_k(x_k, m_k(x_k), x_{k+1}) \mid x_i = x_i]$$

Then the truncated policy:-

$$\pi^{*i} = \{\mu_i^*, \mu_{i+1}^*, \dots, \mu_{N-1}^*\}$$

is optimal for the subproblem.

* Principle of Optimality :-

This is called Principle of Optimality suggests a dynamic programming algo.

Proposition :- For every initial state x_0 , the optimal cost

$J^*(x_0) = J_0(x_0)$ that is given by the last step of the

following algo. that goes backward in time.

$$J_N(x_N) = g_N(x_N). \quad \forall x_N \in S \rightarrow (1)$$

$$J_k(x_k) = \min_{a_k \in A(x_k)} E_{x_{k+1}} [g_k(x_k, a_k, x_{k+1}) + J_{k+1}(x_{k+1})]$$

$\forall x_k \in S$

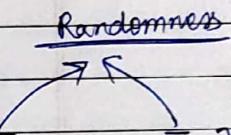
$$\forall k = N-1, N-2, \dots, 0 \rightarrow (2)$$

Also, if $\mu_k^* = \mu_k(x_k)$ minimizes the RHS of (2),
 $\forall k$, the policy $\pi^* = \{\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*\}$

is optimal.

Let's look at (1) & (2) in detail :-

$$J_N(x_N) = g_N(x_N). \quad \forall x_N \in S$$



$$J_{N-1}(x_{N-1}) = \min_{a_{N-1} \in A(x_{N-1})} E_{x_N} [g_{N-1}(x_{N-1}, a_{N-1}, x_N) + J_N(x_N)]$$

* Proof :- Why does this DP algorithm converges.

Let for any admissible policy $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$

$\pi^k = \{\mu_k, \mu_{k+1}, \dots, \mu_{N-1}\}$, $k = 0, 1, \dots, N-1$ denote the

k^{th} tail policy.

Let

$$J_k^*(x_k) = \min_{\pi^k} E_{x_{k+1}, \dots, x_N} [g_k(x_k) + \sum_{i=k}^{N-1} g_i(x_i, \mu_i(x_i))]$$

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

Also, let $J_N^*(x_N) = g_N(x_N) = J_N(x_N) + x_N$.

We show by induction that $J_k^*(x_k) \leq J_R(x_k) \quad \forall k, \forall x_k$.

Assume that for some k and all x_{k+1} , $J_{k+1}^*(x_{k+1}) = J_{k+1}(x_{k+1})$.

Note: $\pi^k = \{\mu_k, \pi^{k+1}\}$, thus $\neq x_k$

$$J_k^*(x_k) = \min_{\{\mu_k, \pi^{k+1}\}} E_{x_{k+1}, \dots, x_N} [g_N(x_N) +$$

$$g_k(x_k, \mu_k(x_k), x_{k+1})$$

$$+ \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1})$$

$$= \min_{\{\mu_k, \pi^{k+1}\}} E_{x_{k+1}} [g_k(x_k, \mu_k(x_k), x_{k+1})]$$

$$+ E_{x_{k+2}, \dots, x_N} [g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1})]$$

$$| x_{k+1} | | x_k = x_k]$$

$$= \min_{\mu_k} E_{x_{k+1}} [g_k(x_k, \mu_k(x_k), x_{k+1}) +$$

$$\min_{\pi^{k+1}} E_{x_{k+2}, \dots, x_N} [g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1})]$$

$$J_{k+1}^*(x_{k+1})$$

$$| x_{k+1}]$$

$$| x_k = x_k]$$

$$= \min_{\mu_k} E_{x_{k+1}} [g_k(x_k, \mu_k(x_k), x_{k+1}), \\ J_{k+1}^*(x_{k+1}) \mid X_k = x_k].$$

from induction hypothesis. $J_{k+1}^*(x_{k+1}) = J_{k+1}(x_{k+1})$.

$$J_k^*(x_k) = \min_{\mu_k} E_{x_{k+1}} [g_k(x_k, \mu(x_k), x_{k+1}) + \\ J_{k+1}(x_{k+1}) \mid X_k = x_k].$$

$$= \min_{q_k \in A(x_k)} E_{x_{k+1}} [g_k(x_k, q_k, x_{k+1}) + \\ J_{k+1}(x_{k+1}) \mid X_k = x_k] \\ = J_k(x_k).$$

$$\Rightarrow J_j^*(x_j) = J_j(x_j) \neq j = 0, 1, \dots, N \neq x_j$$

- i. e. x_j is a minimum point.

$$2) w \neq (w) \text{ if } (w) w = (w) w$$

$$(w) w + (w) w = (w) w$$

$$(w) w = (w) w$$

$$(w) w = (w) w$$

$$(w) w = (w) w$$

Dyn. Prog.
Opt. Control

Finite Horizon Problems :-

In general,

S_k = state space at time k .

A_k = Action space

$A_k(x_k)$ = set of feasible actions at time k in state x_k .

S - State Space.

A - Action Space.

x_k = state at time k .

a_k = action at

$a_k \in \underbrace{A(x_k)}_{\text{Set of feasible}} \subset A$

Set of feasible
actions in x_k

$g_k(x_k, u_k, x_{k+1})$ - simple state cost at $k=0, 1, \dots, N-1$

$g_N(x_N)$ - terminal cost.

* Dynamic Programming Algorithm :-

$$J_N(x_N) = g_N(x_N) \quad \forall x_N \in S.$$

$$J_k(x_k) = \min_{a_k \in A(x_k)} E_{x_{k+1}} [g_k(x_k, a_k, x_{k+1}) + J_{k+1}(x_{k+1})]$$

$\forall k = N-1, \dots, 0$
 $\forall x_k \in S$.

$$= \min_{a_k \in A(x_k)} \sum_{j \in S} p(x_k, a_k, j) (g_k(x_k, a_k, j) + J_{k+1}(j))$$

Analog for infinite horizon.

$$J(x) = \min_{a \in A(x)} \sum_{j=0}^{\infty} p(x, a, j) [g(x, a, j) + \gamma J(j)],$$

$$\gamma \in (0, 1)$$

* Infinite horizon problem (discounted cost).

Suppose π is policy.

$$V_{\pi}(x) = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k g(x_k, \pi(x_k), x_{k+1}) \mid x_0 = x \right]$$

(value of stat x
under policy π)

$$\text{Suppose } |g(x_k, \pi(x_k), x_{k+1})| \leq B$$

$$\Rightarrow |V_{\pi}(x)| \leq \frac{B}{1-\gamma} < \infty$$

* Long run avg costs :-

$$J(\pi) = \lim_{N \rightarrow \infty} \frac{1}{N} E_{\pi} \left[\sum_{k=0}^{N-1} g(x_k, \pi(x_k), x_{k+1}) \mid x_0 = x \right]$$

Exists if under π the markov chains is ergodic.

$$J(\pi) = \lim_{\gamma \rightarrow 1^-} (1-\gamma) V_{\pi}(x)$$

* Example (Finite Horizon) :-

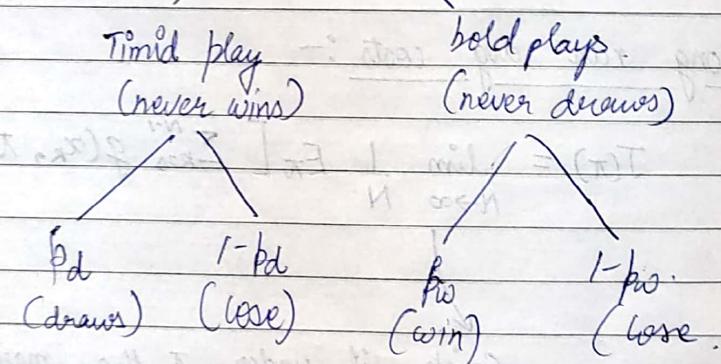
There is a player playing chess with an opponent. Total of N games to be played in match. If scores are tied after N games, then match goes in sudden death mode. (first player after N games to win, wins the match).

Score assignment :-	Win	- 1 pt.
	Loss	- 0 pt.
	Draw	- $\frac{1}{2}$ pt.

State = net score = pts of player - pts of opponent.

Note:- Maximization problem.

Action :- A player can select



Assume :- $r_k(x_k, u_k, x_{k+1}) = 0 \quad \forall k=0, 1, \dots, N-1$

$$r_N(x_N) = \begin{cases} 1 & ; x_N > 0 \\ fw & ; x_N = 0 \\ 0 & ; x_N < 0 \end{cases}$$

Dynamic Programming :- eqn :-

$$J_k(x_k) = \max [p_d \cdot J_{k+1}(x_k) + (1-p_d) J_{k+1}(x_k-1),$$

$$p_w J_{k+1}(x_k+1) + (1-p_w) J_{k+1}(x_k-1)]$$

Thus it is optimal to play bold if,

$$p_w J_{k+1}(x_k+1) + (1-p_w) J_{k+1}(x_k-1).$$

$$\geq p_d J_{k+1}(x_k) + (1-p_d) J_{k+1}(x_k-1).$$

$$p_w \cdot J_{k+1}(x_k+1) + J_{k+1}(x_k-1) - p_w \cdot J_{k+1}(x_k-1)$$

$$\geq p_d \cdot J_{k+1}(x_k) + J_{k+1}(x_k-1) - p_d (J_{k+1}(x_k-1))$$

$$\Rightarrow \frac{p_w}{p_d} \geq \frac{J_{k+1}(x_k) - J_{k+1}(x_k-1)}{J_{k+1}(x_k+1) - J_{k+1}(x_k-1)}$$

Suppose $p_d > p_w$.

Consider $k = N-1$

Consider $k = N-1$

$$J_{N-1}(x_{N-1}) = \max [p_d \cdot J_N(x_{N-1}) + (1-p_d) \cdot J_N(x_{N-1}-1), \\ p_w \cdot J_N(x_{N-1}+1) + (1-p_w) J_N(x_{N-1}-1)]$$

$$J_{N-1}(x_{N-1}) = \begin{cases} 1 & \text{if } x_{N-1} > 1 - \text{Either is optimal} \\ p_d + (1-p_d)p_w & \text{if } x_{N-1} = 1 - \text{Play timid} \\ p_w & \text{if } x_{N-1} = 0 - \text{Play bold} \\ p_w^2 & \text{if } x_{N-1} = -1 - \text{play bold} \\ 0 & \text{if } x_{N-1} < -1 - \text{Play timid} \end{cases}$$

for 2 games remaining :-

$$J_{N-2}(x_{N-2}) = \max \left\{ p_d \cdot J_{N-1}(x_{N-2}) + (1-p_d) \cdot J_{N-1}(x_{N-2}-1) \right. \\ \left. p_w \cdot J_{N-1}(x_{N-2}+1) + (1-p_w) \cdot J_{N-1}(x_{N-2}-5) \right\}$$

$$J_{N-2}(x_{N-2}) = \begin{cases} p_w & \text{if } x_{N-2} = 0 - \text{play bold.} \\ 1-p_d + (1-p_w) & \text{if } x_{N-2} > 0 \rightarrow \text{play timid.} \\ J_{N-1}(x_{N-2}-1). \end{cases}$$

$$(1-p_d) \cdot J_{N-1}(x_{N-2}) + (1-p_w) \cdot J_{N-1}(x_{N-2}-1) \\ (1-p_d) \cdot J_{N-1}(x_{N-2}-1) + (1-p_w) \cdot J_{N-1}(x_{N-2}-5)$$

~~play bold + (1-p_w)~~.

$$\cancel{p_d} + (1-p_d) J_{N-1}(x_{N-2}-1) \\ \cancel{p_w} + (1-p_w) J_{N-1}(x_{N-2}-5)$$

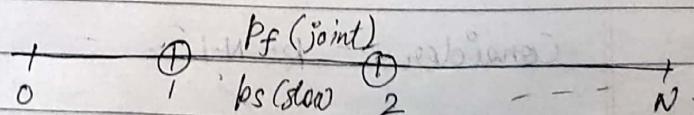
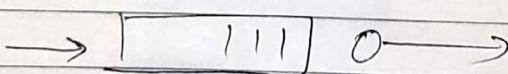
$$(1-p_d) \cdot J_{N-1}(x_{N-2}) = (1-p_d) \cdot J_{N-1}(x_{N-2}) + (1-p_w) \cdot J_{N-1}(x_{N-2}-5)$$

$$n + (1-n)a$$

$$g + (1-g)a$$

$$(1-a)x + \cancel{a} \\ (1-a)y.$$

Ex.



buffer size = n

probability that m customers come in at an instant = p_m .

Customer served as per FCFS:-

Types of service → fast :- cost μ_f - exists w.p. p_f
 $(\mu_f > \mu_s)$

→ slow :- cost μ_s - exists w.p. p_s
 $(p_f > p_s)$

If at end of a period, i customers in system.
holding cost = $\gamma(i)$.

If at the end of the horizon, i customers in system,
then cost $\rightarrow R(i)$.

Assume no. of arrivals in a period independent of no. of arrivals in any other period.

Transition Prob.:-

$$p_{0j}(\mu_f) = p_j = p_{0j}(\mu_s), \quad j = 0, 1, \dots, n-1.$$

$$p_{nn}(\mu_f) = \sum_{j=n}^{\infty} p_j = p_{nn}(\mu_s).$$

for $i > 0$

$$p_{ij}(\mu_f) = 0 \quad \text{if } j < i-1.$$

$$p_{ij}(\mu_f) = p_f p_0 \quad \text{if } j = i-1$$

\uparrow
Fast
Slow
no. of arrival of people.

$$\begin{aligned} p_{ij}(\mu_f) &= P(i-j+1 \text{ arrivals, service completion}) \quad i-1 < j < n-1 \\ &\quad + P(j-i \text{ arrivals, no service completion}) \\ &= p_f p_{j-i+1} + (1-p_f) p_{j-i} \end{aligned}$$

$j = n-1$

$$p_{ij}^o(\mu_f) = p_f \sum_{m=n-i}^{\infty} p_m + (1-p_f) p_{n-1-i}$$

$j = n$

$$p_{ij}^o(\mu_f) = (1-p_f) \sum_{m=n-i}^{\infty} p_m$$

DP Algo :-

$$J_N(i) = R(i), \quad i = 0, 1, \dots, n.$$

$$J_K(x_k) = \min \left[r(i) + c_f + \sum_{j=0}^{\infty} p_{ij}^o(\mu_f) J_{K+1}(j), \right.$$

$$\left. r(i) + c_s + \sum_{j=0}^{\infty} p_{ij}^o(\mu_s) J_{K+1}(j) \right]$$

(2)

$a < i < d$

$i < d$

$a = (II) \text{ if}$

$i = d$

$d = (II) \text{ if}$

$\max\{x_{i-1}, \text{ (minimum value of } i+1 \text{ to } d) \} = (II) \text{ if}$

$p_{ii}(d-1) + \min_{i \leq j \leq d-1}$



Scanned with OKEN Scanner