

Finite horizon Markov decision processes (MDPs)

(Lecture - I)

Example 1: Machine replacement

A problem over N stages
(e.g. think of maintaining a bus, with each stage corresponding to a month).

The machine can be in one of the " n " states, i.e.,

$\{1, \dots, n\}$
↑ perfect condition ↑ worst condition

"Don't confuse state with stage"

Operating cost: $g(i)$ in state i

$$g(1) \leq g(2) \leq \dots \leq g(n)$$

Actions:

"Repair"

↓
machine becomes new (i.e., goes to state 1)

or "Do nothing"

"machine gets progressively worse"

"Stochastic system"

P_{ij} : probability that the machine transitions from state i to j "when you do nothing".

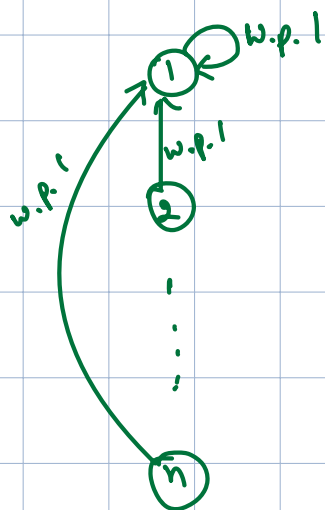
→ transition probabilities

$$\sum_{j=1}^n P_{ij} = 1, \quad P_{ij} = 0 \text{ if } j < i$$

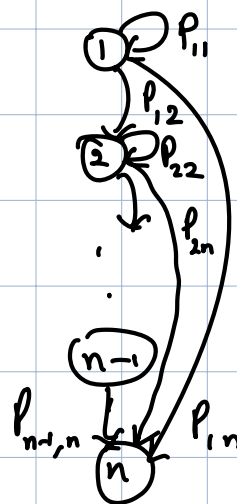
If you repair, the machine transitions to state "1" & stays there for one stage

Transition diagram:

"Repair"



"Do nothing"

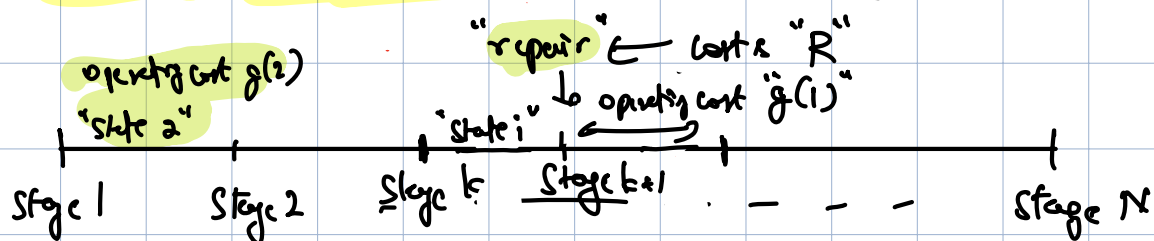


w.p.
= with
Probability

$$p_{ij} = 0 \quad \text{if } j < i$$

On "repair", machine goes to state "1", & remains there for one stage. Repair cost is "R".

Goal: choose actions so that the total cost (i.e., cumulative cost over N stages) is minimized.

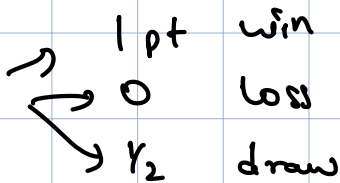


Takeaway: MDP has "states", "actions", "transition probabilities", & "costs".

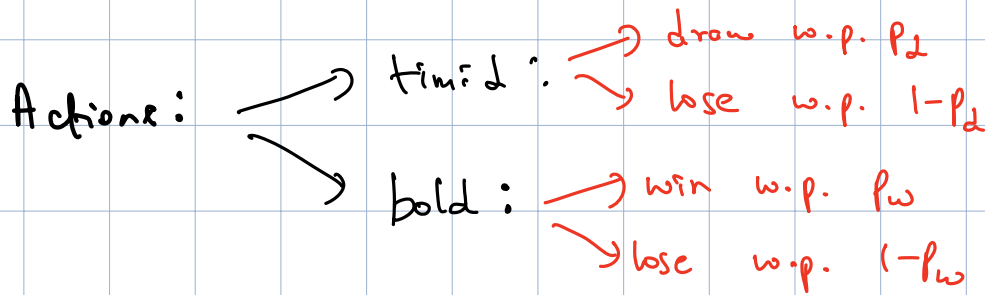
Another example:

"Chess match"

"Fix opponent"
You play 2 games

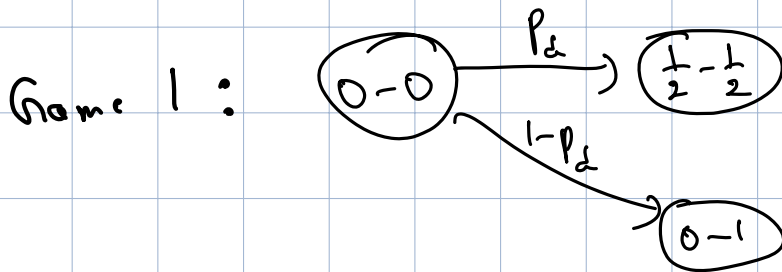


If there is a tie, then
"sudden death" phase where
you play games one after the other
until a decisive result.

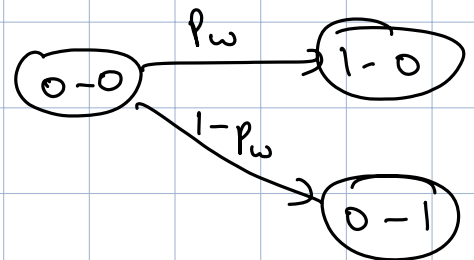


" $p_d > p_w$ "

"Timid play"

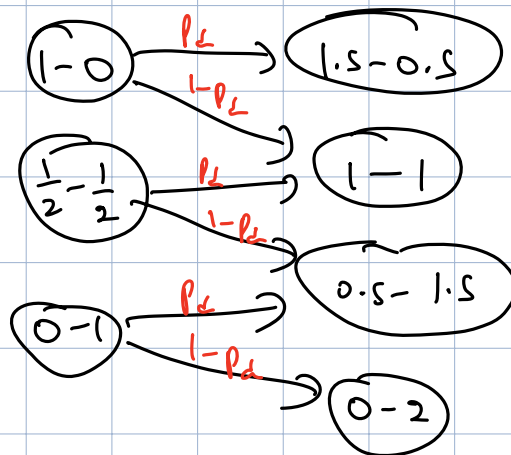


"Bold play"

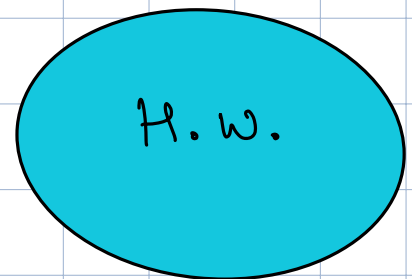


"Timid"

Game 2:



"Bold"



Sudden Death: Play bold.

MDP framework:

Let \mathcal{X} denote the state space,
 \mathcal{A} denote the action space.

$x_k \rightarrow$ state in stage k , $x_k \in \mathcal{X}$

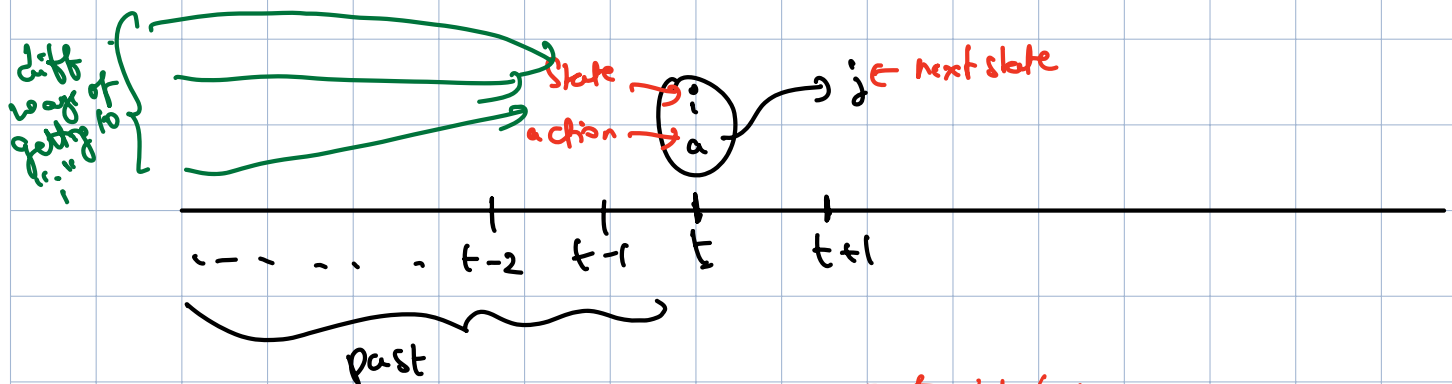
set of possible
actions depend on
the current state

$a_k \rightarrow$ action that can be taken in state x_k , $a_k \in \mathcal{A}(x_k) \subset \mathcal{A}$.

Transition probabilities:

$$p_{ij}(a) = P(x_{k+1} = j \mid x_k = i, a_k = a)$$

State evolution satisfies "Controlled Markov property".



$$P(\overset{\text{future}}{x_{t+1} = j} \mid \overset{\text{current state \& action}}{x_t = i, a_t = a}, \overset{\text{past states/actions}}{x_{t-1}, a_{t-1}, \dots, x_0})$$

$$= p_{ij}(a)$$

"Single-stage cost" $g_k(i, a, j)$

In stage k , if you are in state i & chose to take action a to transition to state j , then cost incurred is $g_k(i, a, j)$

$g_N(i) \rightarrow$ Terminal cost
 final stage if state is i ,
 then cost is $g_N(i)$

Policy

$\{\mu_0, \dots, \mu_{N-1}\}$ state $\rightarrow \boxed{\mu_k} \rightarrow$ action
 $\mu_k: \mathcal{X} \rightarrow \mathcal{A}$ is a mapping from states to actions for stage k .

Lecture-2*

"Admissible policy": $\mu_k(i) \in \mathcal{A}(i) \leftarrow$ set of feasible action in state i

We say a policy $\pi = \{\mu_0, \dots, \mu_{N-1}\}$ is admissible if $\mu_k(i) \in \mathcal{A}_k(i) \quad \forall k=0, \dots, N-1 \text{ \& } \forall i \in \mathcal{X}$.

Total cost: Initial state $x_0 \in \mathcal{X}$

$$J_\pi(x_0) = E_{x_1, \dots, x_N} \left[\underbrace{g_N(x_N)}_{\substack{\downarrow \\ \text{terminal cost}}} + \sum_{k=0}^{N-1} \underbrace{g_k(x_k, \mu_k(x_k), x_{k+1})}_{\substack{\text{cost in stage } k}} \right]$$

$\begin{matrix} \text{current state} & \downarrow & \text{action chosen} & \downarrow & \text{next state} \\ & & & & \end{matrix}$

$\xrightarrow{\text{state transition here}}$
 Stage k Stage $k+1$

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

"Optimization objective"

optimal total cost

$$J_{\pi^*}(x_0) =$$

\uparrow
 optimal policy

$$\min_{\pi \in \Pi}$$

\downarrow

set of "admissible" policies

$$J_\pi(x_0),$$

total cost of policy π

$$\forall x_0 \in \mathcal{X}$$

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{argmin}} J_\pi(x_0)$$

Open vs. closed loop policies:

Closed loop policy: $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$

μ_i : decision based on the state
 $\mu_i(x_i)$, $x_i \in \mathcal{X}$.

Open loop policy: fix the sequence of actions beforehand

Chess match (revisited):

Possible open loop policies:

	Game 1	Game 2	Prob(win)?
Policy π_1	Timid	Timid	$p_d^2 p_w$
Policy π_2	Bold	Bold	$p_w^2 + 2p_w^2(1-p_w)$
Policy π_3	Bold	Timid	$p_w p_d + p_w^2(1-p_d)$
Policy π_4	Timid	Bold	$p_w p_d + p_w^2(1-p_d)$

Let's ignore π_1 (If $3p_w > p_d$, then π_2 is better than π_1)

Which among π_2, π_3, π_4 is the best?

$$\begin{aligned} & \max(p_w^2(3-2p_w), p_w p_d + p_w^2(1-p_d)) \\ &= \max(p_w^2 + 2p_w^2(1-p_w), p_w^2 + p_w p_d(1-p_w)) \\ &= p_w^2 + p_w(1-p_w) \max(2p_w, p_d) \end{aligned}$$

If $2p_w < p_d$, then π_3/π_4 are better

Else, π_2 is better

Set $p_w = 0.45$, $p_d = 0.9$, Then $\text{Prob}(\text{winning match}) = 0.425$
 $< 50\%$ chance of match win

Closed loop policy: Play timid if leading, else play bold.
 π_c

$$\begin{aligned} & \text{Prob}(\text{match win with } \pi_c) \\ &= p_w p_d + p_w (p_w (1-p_d) + (1-p_w) p_w) \\ &= p_w^2 (2-p_d) + p_w (1-p_w) p_d \end{aligned}$$

If $p_w = 0.45$, $p_d = 0.9$, then $\text{Prob}(\text{winning match with } \pi_c)$
 $= 0.53 \rightarrow > 50\%$ chance of winning match

Optimality principle

Let $\pi^* = \{\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*\}$ denote the optimal policy.

Consider the tail sub-problem

$$J_i^*(x_i) = \min_{\pi^i = (\mu_i, \dots, \mu_{N-1})} E \left(g_N(x_N) + \sum_{k=i}^{N-1} g_k(x_k, \mu_k(x_k), x_{k+1}) \right)$$

Optimality principle

The policy $\{\mu_i^*, \dots, \mu_{N-1}^*\}$ is optimal for the
(N-i) stage problem in (*)

Proof: see Sec 1.5 of
Bertsekas DPOC Vol. I

DP algorithm:

$$\text{Set } J_N(x_N) = g_N(x_N), \forall x_N \in X$$

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

$$\left\{ \begin{array}{l} J_k(x_k) = \min_{a_k \in A(x_k)} E_{x_{k+1}} \left(g_k(x_k, a_k, x_{k+1}) + J_{k+1}(x_{k+1}) \right) \\ \end{array} \right. \quad \forall x_k \in X.$$

solⁿ of tail problem

Idea: Going backward, J_0 from DP algorithm is the optimal cost i.e., $J_{\pi^*} \equiv J_0$

Applying DP algorithm to "machine replacement" example:

$$J_N(1) = 0 \quad (\text{no terminal cost})$$

for $i=1, \dots, n$, $J_k(i) = \min \left(\underbrace{R + g(1) + J_{k+1}(1)}_{\text{repair}}, \underbrace{g(i) + \sum_{j=i}^n P_{ij} J_{k+1}(j)}_{\text{do nothing}} \right)$

Lecture-3

"DP algorithm finds the best policy"

Claim: $\forall x_0 \in X$, the function $J_0(x_0)$ obtained at the end of the DP algorithm coincides with the optimal cost $J^*(x_0) (= J_{\pi^*}(x_0))$.

Proof: For any admissible policy $\pi = \{\mu_0, \dots, \mu_{N-1}\}$,
let $\pi^k = \{\mu_k, \dots, \mu_{N-1}\}$

$J_k^*(x_k)$ be the optimal cost of the tail sub-problem beginning in stage k , in state x_k .

$$J_k^*(x_k) = \min_{\pi^k = \{\mu_{k+1}, \dots, \mu_{N-1}\}} E \left(g_N(x_N) + \sum_{i=k}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1}) \right)$$

(claim: $J_k^*(x_k) = \mathcal{J}_k^*(x_k)$, $\forall k$)
 ↳ obtained by DP algorithm

< Pf when Pf >

Base case: $J_N^*(x_N) = g_N(x_N) = \mathcal{J}_N(x_N)$, $\forall x_N \in X$

Induction hypothesis: Assume the claim for " $k+1$ ", i.e.,
 $J_{k+1}^*(x_{k+1}) = \mathcal{J}_{k+1}^*(x_{k+1})$, $\forall x_{k+1}$

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

$$= \min_{\mu_k} E_{x_{k+1}} \left(g_k(x_k, \mu_k(x_k), x_{k+1}) + \min_{\pi^{k+1}} E_{x_{k+1}, \dots, x_N} \left[g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1}) \right] \middle| x_{k+1} \right)$$

splitting of mins \rightarrow optimality principle (for - rigorous proof, check Sec 1.5 of Dpoc Vol. I)

$$J_k^*(x_k) = \min_{\mu_k} E_{x_{k+1}} \left(g_k(x_k, \mu_k(x_k), x_{k+1}) + J_{k+1}^*(x_{k+1}) \right)$$

Induction hypothesis
 \downarrow

$$= \min_{\mu_k} E \left(g_k(x_k, \mu_k(x_k), x_{k+1}) + J_{k+1}(x_{k+1}) \right)$$

$$= \min_{a_k \in \mathcal{A}(x_k)} E \left(g_k(x_k, a_k, x_{k+1}) + J_{k+1}(x_{k+1}) \right)$$

$$= J_k^*(x_k)$$

To infer this, we used

$$\min_{\mu \in B} F(x, \mu(x)) = \min_{a \in \mathcal{A}(x)} F(x, a)$$

$$B = \{ \mu \mid \mu(x) \in \mathcal{A}(x) \}$$

$$\text{So, } J_k^*(x_k) = J_k(x_k)$$

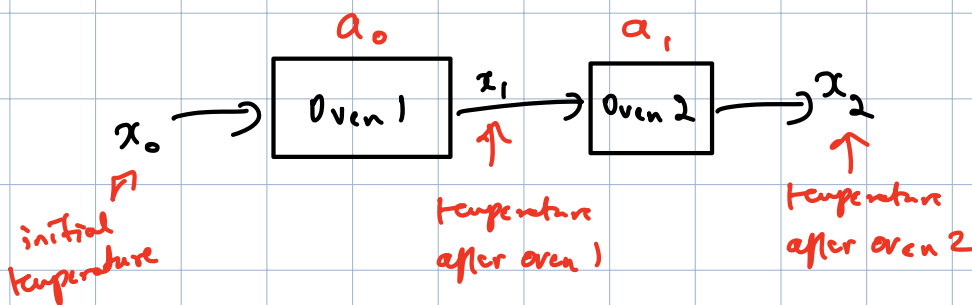
$$\Rightarrow J_0(x_0) = J^*(x_0)$$

End of pf within pf ■

Examples:

Lecture-4*

① Linear system with quadratic cost



Goal: to get x_2 as close as possible to a target temperature T

State evolution :

$$x_{k+1} = (1-\alpha)x_k + \alpha a_k, \quad k = 0, 1, \dots, (\infty)$$

linear state evolution

$\alpha \in (0, 1)$
fixed

Total cost : $a_0^2 + a_1^2 + (x_2 - T)^2$
(to be minimized)

Apply DP algorithm:

Final stage:

$$J_2(x_2) = (x_2 - T)^2 \quad \text{Terminal cost}$$

Going back :
one stage

$$J_1(x_1) = \min_{a_1} (a_1^2 + J_2(x_2))$$

$$= \min_{a_1} (a_1^2 + J_2((1-\alpha)x_1 + \alpha a_1))$$

$$= \min_{a_1} (a_1^2 + ((1-\alpha)x_1 + \alpha a_1 - T)^2)$$

Linear in x_1

$$\mu_1^*(x_1) = \frac{\alpha (T - (1-\alpha)x_1)}{1 + \alpha^2}$$

$$J_1^*(x_1) = \frac{((1-\alpha)x_1 - T)^2}{1 + \alpha^2}$$

→ quadratic in x_1

Check:

$$\mu_0^*(x_0) = \frac{(1-\alpha)\alpha(T - (1-\alpha)^2 x_0)}{1 + \alpha^2(1 + (1-\alpha)^2)} \quad \leftarrow \text{linear in } x_0$$

$$J_0^*(x_0) = \frac{((1-\alpha)^2 x_0 - T)^2}{1 + \alpha^2(1 + (1-\alpha)^2)} \quad \leftarrow \text{quadratic in } x_0$$

Adding randomness to ovens:

Stochastic
"shut"
evolution

$$\rightarrow x_{k+1} = (1-\alpha)x_k + \alpha a_k + \omega_k, \quad k=0,1$$

↑
zero-mean r.v. iid (indep of $\{x_k\}$)
with bounded variance
(e.g. $N(0, \sigma^2)$)

Applying DP algorithm:

$$\begin{aligned} J_1(x_1) &= \min_{a_1} E_{\omega_1} \left(a_1^2 + ((1-\alpha)x_1 + \alpha a_1 + \omega_1 - T)^2 \right) \\ &= \min_{a_1} \left(a_1^2 + ((1-\alpha)x_1 + \alpha a_1 - T)^2 \right. \\ &\quad \left. + 2 E \omega_1 ((1-\alpha)x_1 + \alpha a_1 - T) \right. \\ &\quad \left. + E \omega_1^2 \right) \\ &= \min_{a_1} \left(a_1^2 + ((1-\alpha)x_1 + \alpha a_1 - T)^2 + E \omega_1^2 \right) \end{aligned}$$

Minimizing RHS above leads to the same action as in the deterministic setting, i.e.,

$$\mu_1^*(x_1) = \frac{\alpha(T - (1-\alpha)x_1)}{1+\alpha^2}$$

(hus match - revisited (last time))

Consider an extension to N games

timid $\rightarrow p_L$, bold $\rightarrow p_W$

$$p_L > p_W$$

Players play N games &
enter sudden death if the score is tied.

State: net score (e.g. (0-1) state = -1)

Apply DP algorithm:

(*)
$$J_k(x_k) = \max \left(\underbrace{p_d J_{k+1}(x_k) + (1-p_d) J_{k+1}(x_k-1)}_{\text{timid play}}, \underbrace{p_w J_{k+1}(x_k+1) + (1-p_w) J_{k+1}(x_k-1)}_{\text{bold play}} \right)$$

*we are playing for rewards
(choose min to max in DP algo)*

$$J_N(x_N) = \begin{cases} 1 & \text{if } x_N > 0 \\ p_w & \text{if } x_N = 0 \\ 0 & \text{if } x_N < 0 \end{cases}$$

It is better to play bold when

$$\frac{p_w}{p_d} > \frac{J_{k+1}(x_k) - J_{k+1}(x_k-1)}{J_{k+1}(x_k+1) - J_{k+1}(x_k-1)} \quad \rightarrow \text{inferred from (*)}$$

Given that we have J_N specified, we can go back & calculate J_{N-1} .

x_{N-1}	T_{N-1}	Best play
> 1	1	does not matter
1	$\max(\overbrace{p_d + (1-p_d)p_w}^{\text{Timid}}, \overbrace{p_w + (1-p_w)p_w}^{\text{Bold}})$ $= p_d + (1-p_d)p_w$	Timid
0	p_w	Bold
-1	p_w^2	Bold
< -1	0	does not matter

For the 2-game match \rightarrow we can figure the optimal strategy by knowing $T_{N-2}(0)$

$$\begin{aligned}
 T_{N-2}(0) &= \max\left(\underbrace{p_d p_w + (1-p_d)p_w^2}_{\text{Timid}}, \underbrace{p_w(p_d + (1-p_d)p_w + (1-p_w)p_w^2)}_{\text{bold play}}\right) \\
 &= \max\left(p_w(p_d + (1-p_d)p_w), p_w(p_d + (1-p_d)p_w + (1-p_w)p_w^2)\right) \\
 &= p_w(p_d + (1-p_d)p_w + (1-p_w)p_w^2) \Rightarrow \text{play bold}
 \end{aligned}$$

As noted before, one could choose $p_w < 0.5$ & still get a better than 50-50 chance of winning the match if $p_w(p_d + (1-p_d)p_w + (1-p_w)p_w^2) > 0.5$

Another example: (Job scheduling)

N jobs to schedule

$T_i \rightarrow$ time taken for i th job to complete

T_i is a r.v. $\{T_i, i=1, \dots, N\}$ independent

Each job " i " has a reward R_i associated with it.

So, if job " i " finishes at time " t ", then

the reward is $\alpha^t R_i$, $\alpha =$ discount factor $0 < \alpha < 1$

Cumulative reward = sum of each job's reward.

Goal: schedule jobs to maximize cumulative reward.

"Interchange argument" to figure optimal schedule

$$L = \{i_0, i_1, \dots, i_{k-1}, i, j, i_{k+2}, \dots, i_{N-1}\}$$

$$L' = \{i_0, i_1, \dots, i_{k-1}, j, i, i_{k+2}, \dots, i_{N-1}\}$$

$$J_L = E \left[\alpha^{t_0} R_{i_0} + \dots + \alpha^{t_{k-1}} R_{i_{k-1}} + \alpha^{t_{k-1} + T_i} R_i + \alpha^{t_{k-1} + T_i + T_j} R_j + \dots + \alpha^{t_{N-1}} R_{i_{N-1}} \right]$$

$$J_{L'} = E \left[\alpha^{t_0} R_{i_0} + \dots + \alpha^{t_{k-1}} R_{i_{k-1}} + \alpha^{t_{k-1} + T_j} R_j + \alpha^{t_{k-1} + T_j + T_i} R_i + \dots + \alpha^{t_{N-1}} R_{i_{N-1}} \right]$$

Schedule L is better than L' if

$$E \left[\alpha^{t_{k-1} + T_i} R_i + \alpha^{t_{k-1} + T_i + T_j} R_j \right] \geq E \left[\alpha^{t_{k-1} + T_j} R_j + \alpha^{t_{k-1} + T_j + T_i} R_i \right]$$

Using t_{k-1}, T_i, T_j are independent,

$$\frac{E(\alpha^{T_i}) R_i}{1 - E(\alpha^{T_i})} \geq \frac{E(\alpha^{T_j}) R_j}{1 - E(\alpha^{T_j})} \quad (*)$$

From (*), the optimal schedule works out as follows:

Assign $\mu_i = \frac{E(\alpha^{T_i}) R_i}{1 - E(\alpha^{T_i})}$ as the index for job i , $i=1 \dots N$

Order $\{\mu_1, \dots, \mu_N\}$, say $\mu_{[1]} \geq \mu_{[2]} \dots \geq \mu_{[N]}$

Optimal schedule = $\{[1], [2], \dots, [N]\}$

index-based optimal policy

Further reading: Check out
Gittin's index
Sec 1.5 of PLOC v.1.1

Yet-another example: <Optimal stopping>
"Asset-selling"

A technical note before asset-selling:

Discrete-time MDPs can be formulated as

$$x_{k+1} = f(x_k, a_k, \omega_k)$$

↪ disturbance

$\{w_k\}$ could be i.i.d or could depend on x_k, a_k
 $x_k \in$ infinite set.

For the case when $x_k \in \{1, \dots, n\}$, it is enough to know $p_{ij}^a = P(x_{k+1}=j | x_k=i, a_k=a)$

Now to asset-selling!

Want to sell an asset.

You get offers w_0, w_1, \dots, w_{N-1}

Assume: $\{w_k\}$ iid with some distribution function that has finite mean

Action \rightarrow sell the asset (by accepting the offer) a^1
 \rightarrow don't sell & wait for more offers a^2

Add a special state "T" to denote that the asset is sold.

Also, $x_0 = 0$

$x_{k+1} = f(x_k, a_k, w_k)$, where

$$f(x_k, a_k, w_k) = \begin{cases} T & \text{if } (x_k \neq T, a_k = a^1) \text{ or } (x_k = T) \\ w_k & \text{else.} \end{cases}$$

Work with rewards.

Goal: maximize $E \left(g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, a_k, w_k) \right)$

\nearrow Expectation over w_0, w_1, \dots, w_{N-1}

$$g_N(x_N) = \begin{cases} x_N & \text{if } x_N \neq T \\ 0 & \text{else} \end{cases}$$

$$g_k(x_k, a_k, w_k) = \begin{cases} (1+r)^{N-k} x_k & \\ 0 & \end{cases}$$

"interest added"

asset not sold *sell*

$$x_k \neq T, a_k = a'$$

else

$1 > r > 0$ is the interest rate.

Apply DP algorithm:

$$J_N(x_N) = \begin{cases} x_N & \text{if } x_N \neq T \\ 0 & \text{else} \end{cases}$$

$$J_k(x_k) = \begin{cases} \max\left(\underbrace{(1+r)^{N-k} x_k}_{\text{sell}}, \underbrace{E(J_{k+1}(w_k))}_{\text{don't sell}}\right) & \text{if } x_k \neq T \\ 0 & \text{if } x_k = T \end{cases}$$

Then is x_{k+1} if one doesn't sell

$$\text{Let } \alpha_k = \frac{E(J_{k+1}(w_k))}{(1+r)^{N-k}} = \frac{E(J_{k+1}(w))}{(1+r)^{N-k}}$$

(Since w_k iid)

Optimal policy: threshold-based policy "non-stationary thresholds α_k "

sell	if $x_k > \alpha_k$		if $x_k = \alpha_k$
don't sell	if $x_k < \alpha_k$		both actions are fine

Understanding the optimal policy:

Lecture 5*

Suppose $d_k \geq d_{k+1} \quad \forall k$

For notational convenience, let $V_k(x_k) = \frac{J_k(x_k)}{(1+r)^{N-k}}$,
for $x_k \neq T$

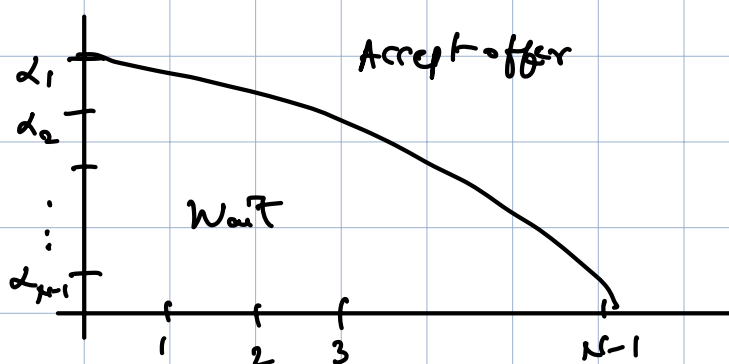
DP algo in "V" notation is

$$V_N(x_N) = x_N$$

$$V_k(x_k) = \max \left(x_k, \frac{E(V_{k+1}(w_k))}{(1+r)} \right)$$

Optimal policy $\mu_k^*(x_k) = \begin{cases} a^1 & \text{if } x_k \geq d_k \\ a^2 & \text{else} \end{cases}$

Claim: $d_k \geq d_{k+1}$



$$d_k = \frac{E(V_{k+1}(w))}{1+r}$$

To show $\alpha_k \geq \alpha_{k+1}$, it is enough if we establish that $V_k(x) \geq V_{k+1}(x) \quad \forall x$

For $k = N-1$,

$$\begin{aligned} V_{N-1}(x) &= \max \left(x, \frac{E(V_N(\omega))}{1+r} \right) \\ &= \max \left(x, \frac{E(\omega)}{1+r} \right) \geq x = V_N(x) \end{aligned}$$

For $k = N-2$,

$$\begin{aligned} V_{N-2}(x) &= \max \left(x, \frac{E(V_{N-1}(\omega))}{1+r} \right) \\ &\geq \max \left(x, \frac{E(V_N(\omega))}{1+r} \right) \\ &= V_{N-1}(x) \end{aligned}$$

Proceeding similarly, we get $V_k(x) \geq V_{k+1}(x), \forall x, \forall k$

Understanding the asset selling problem for large N :

Suppose " ω " is a continuous, positive-valued r.v. with distribution F_ω & density h .

$$V_{k+1}(\omega) = \begin{cases} \omega & \text{if } \alpha_{k+1} \leq \omega \\ \alpha_{k+1} & \text{else} \end{cases} \quad \left. \vphantom{\begin{cases} \omega \\ \alpha_{k+1} \end{cases}} \right\} \text{ rewriting the max in def of } V_{k+1}$$

$$\begin{aligned}
\alpha_k &= \frac{E(V_{k+1}(\omega))}{1+r} \\
&= \frac{1}{1+r} \int_0^\infty V_{k+1}(\omega) h(\omega) d\omega \\
&= \frac{1}{1+r} \int_0^{\alpha_{k+1}} \alpha_{k+1} h(\omega) d\omega + \frac{1}{1+r} \int_{\alpha_{k+1}}^\infty \omega h(\omega) d\omega \\
\alpha_k &= \frac{\alpha_{k+1}}{1+r} F_\omega(\alpha_{k+1}) + \frac{1}{1+r} \int_{\alpha_{k+1}}^\infty \omega h(\omega) d\omega
\end{aligned}$$

$$\alpha_k \leq \frac{\alpha_{k+1}}{1+r} + \frac{E(\omega)}{1+r}$$

used $\int_{\alpha_{k+1}}^\infty \omega h(\omega) d\omega \leq E(\omega)$

$$\alpha_k \stackrel{\alpha_k \geq \alpha_{k+1}}{\leq} \frac{\alpha_k}{1+r} + \frac{E(\omega)}{1+r}$$

From the above, we can conclude

$$0 \leq \alpha_k \leq \frac{E\omega}{r} \quad (*)$$

The sequence $\{\alpha_k\}$ is non-increasing and bounded because of (*).

Since $\alpha_{k+1} \leq \alpha_k$, the sequence $\{\alpha_k\}$ converges as $k \rightarrow \infty$, say to some

$$\bar{\alpha} \in \left[0, \frac{E\omega}{r}\right]$$

$$x_k = \frac{x_{k+1}}{1+r} F_w(x_{k+1}) + \frac{1}{1+r} \int_{x_{k+1}}^{\infty} \omega h(\omega) d\omega$$

Taking $k \rightarrow -\infty$ in eqn above, we obtain

$$\bar{x} = \frac{\bar{x} F_w(\bar{x})}{1+r} + \frac{1}{1+r} \int_{\bar{x}}^{\infty} \omega h(\omega) d\omega$$

Given the distribution F_w of ω , we can calculate \bar{x}

For very large N , an approximately optimal policy is

sell at x_k if $x_k > \bar{x}$
 don't sell at x_k if $x_k \leq \bar{x}$

Simpler "threshold" based policy

because the threshold \bar{x} is the same for all stages.

< End of Finite horizon MDPs >