

# Finite horizon Markov decision processes (MDPs)

## (Lecture - 1\*)

### 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.,

$$\{1, \dots, n\}$$

$\nearrow$   
perfect condition  
 $\uparrow$   
Worst Condition

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

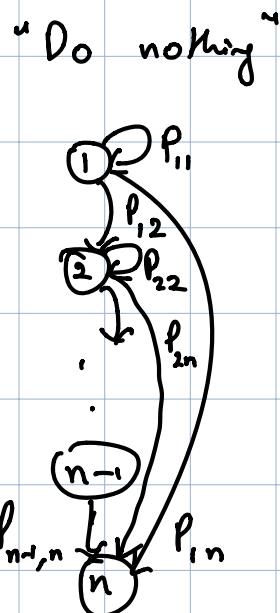
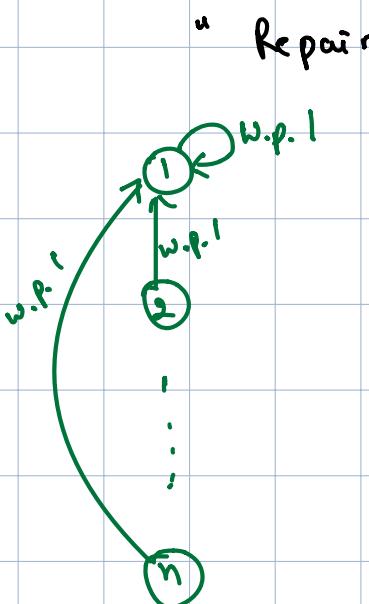
$$g(1) \leq g(2) \leq \dots \leq g(n)$$

Actions: "Repair" or "Do nothing"  
"Repair"  $\rightarrow$  machine becomes new  
"Do nothing"  $\rightarrow$  machine becomes progressively worse

"Stochastic system"

$P_{ij}$ : probability that the machine goes from state  $i$  to state  $j$

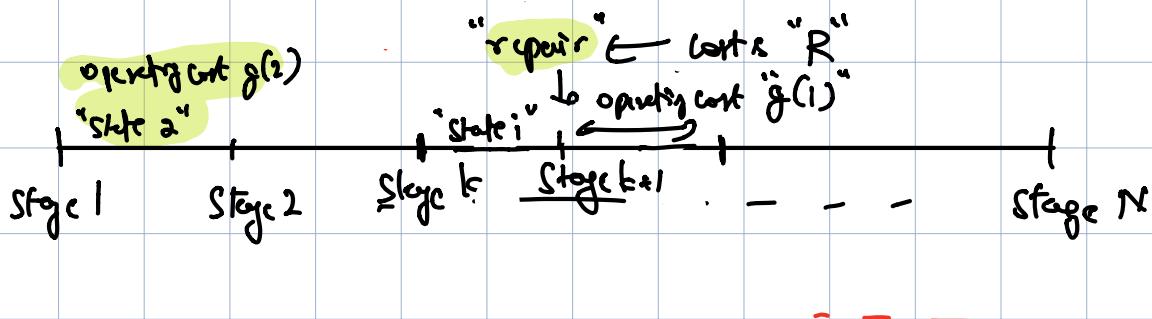
Transition diagram:



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

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

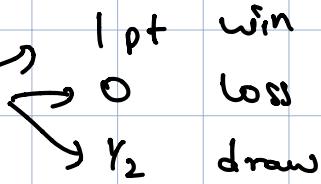
Goal: choose actions so that the total cost (i.e., cumulative cost over N stages) is minimized.



**Takeaway:** MDP has "states", "actions", "stochastic transitions" & the goal is to minimize "total cost".

Another example:

"Fix opponent"  
You play 2 games



"Chess match"

If there is a tie, then

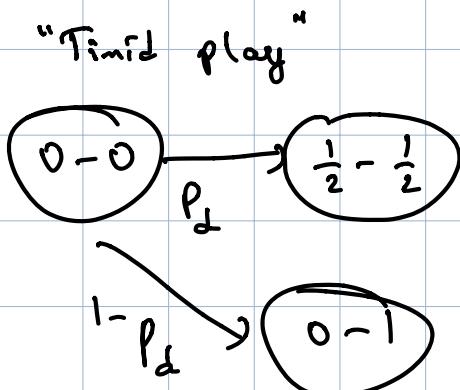
"sudden death" phase where you play games one after the other until a decisive result.

Actions: playing styles

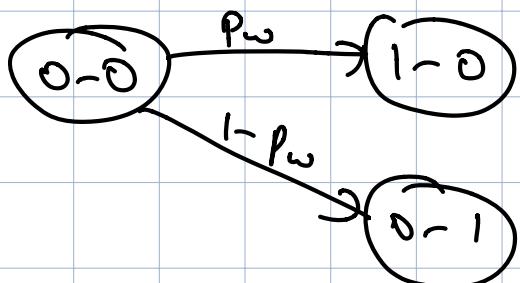
→ timid: draw w.p.  $p_d$  & lose w.p.  $1-p_d$   
→ bold: win w.p.  $p_w$  & lose w.p.  $1-p_w$

" $p_d > p_w$ "

Game 1:

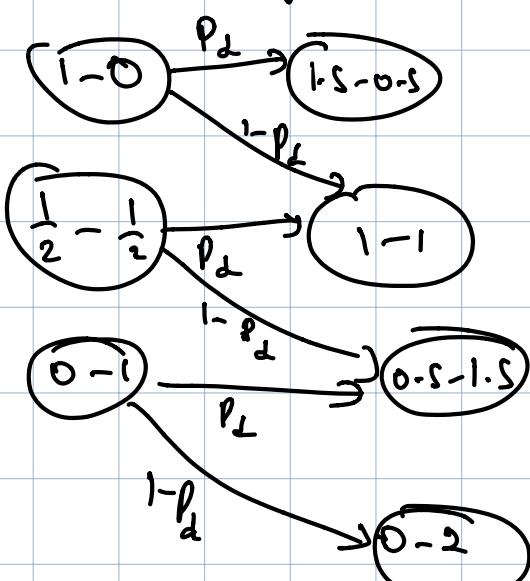


"Bold" play

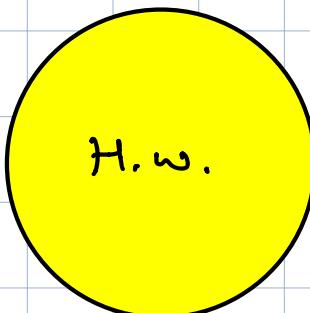


Game 2:

"timid play"



"Bold play"



## MDP framework:

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

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

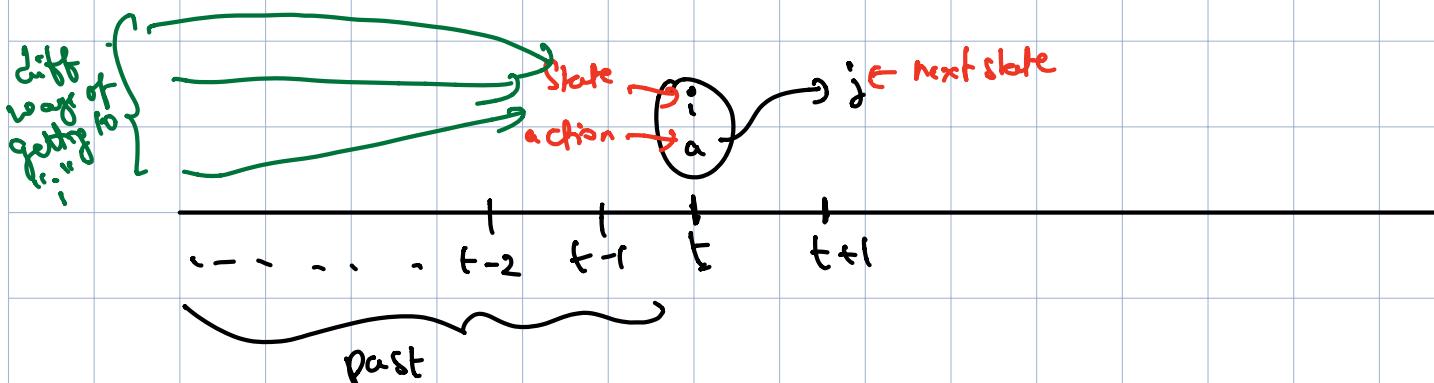
set of possible actions depend on the current state

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

Transition probability:

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

State evolution satisfies "Controlled Markov property".



$$P(x_{t+1} = j \mid x_t = i, a_t = a, x_{t-1}, a_{t-1}, \dots, x_0)$$

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

"Single-stage cost"

$g_k(i, a, j) \leftarrow$  in stage  $k$ , if state  $i$  & action taken is  $a$ , & the next state is  $j$

$g_N(i) \leftarrow$  termination cost in the final stage in state  $i$ ,  
 $i \in S$

**Policy**  $\{ \mu_0, \dots, \mu_{N-1} \}$

$\mu_k : S \rightarrow A$  mapping that specifies how actions are taken in stage  $k$ .

"Admissible policy":  $\mu_k(i) \in A(i)$  [set of feasible actions in state  $i$ ]

$\pi = \{ \mu_0, \dots, \mu_{N-1} \}$  is admissible if  
 $\mu_k(i) \in A(i) \quad \forall k = 0 \dots N-1$

**Total cost**: Initial state  $x_0 \in S$

$$J_\pi(x_0) = \underset{x_1, \dots, x_N}{E} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), x_{k+1}) \right]$$

↓ action chosen  
 ↓ current state  
 ↓ next state  
 ↓ cost in stage  $k$   
 ↓ terminal cost

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

Side note: Fixing  $\pi$ , the sequence of states  $x_0, x_1, \dots, x_N$  is a Markov chain

"Optimization objective"  
optimal total cost  $\rightarrow$  total cost of policy  $\pi$

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

↑ optimal policy  
 ↓ set of "admissible" policies

$$\pi^* = \arg \min_{\pi \in \Pi} J_\pi(x_0)$$

## Lecture - 2\*

**Open vs. closed loop policies:**

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

$\mu_i$ : 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 $\pi_1$	Timid	Timid	$p_d^2 p_w$
Policy $\pi_2$	Bold	Bold	$p_w^2 + 2p_w^2(1-p_w)$
Policy $\pi_3$	Bold	Timid	$p_w p_d + p_w^2(1-p_d)$
Policy $\pi_4$	Timid	Bold	$p_w p_d + p_w^2(1-p_d)$

Let's ignore  $\pi_1$  ( $3p_w > p_d$ , then  $\pi_2$  is better than  $\pi_1$ )

Which among  $\pi_2, \pi_3, \pi_4$  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  $\pi_3/\pi_4$  are better  
else,  $\pi_2$  is better.

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

Closed loop policy: Play timid if leading, else play bold.  
 $\pi_c$

$$\begin{aligned} & \text{Prob( match win with } \pi_c) \\ &= P_w P_d + P_w (P_w(1-P_d) + (1-P_w)P_d) \\ &= P_w^2 (2 - P_w) + P_w (1 - P_w) P_d \\ &= 0.53 \quad \text{for } P_w = 0.45 \text{ and } P_d = 0.9 \\ &> 50\% \text{ chance of match win} \end{aligned}$$

### Optimality principle

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

Consider the tail sub-problem

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

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:

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

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

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

Applying DP algorithm to "machine replacement" example:

$$J_N(i) = 0 \quad (\text{no terminal opt})$$

$$J_k(i) = \min \left( R + g(i) + J_{k+1}(i), g(i) + \sum_{j=i}^n p_{ij} J_{k+1}(j) \right)$$

repair                            do nothing

"DP algorithm finds the best policy"

Claim:  $\forall x_0 \in \mathcal{X}$ , 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  $\pi = \{\mu_0, \dots, \mu_{N-1}\}$ ,

$$\text{let } \pi^k = \{\mu_k, \dots, \mu_{N-1}\}$$

$\mathcal{J}_k^*(x_k)$  be the optimal cost of the tail subproblem beginning in stage  $k$ , in state  $x_k$ .

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

$\pi^k = \{\mu_{k,i} - \mu_{N-1,i}\}$

(claim:  $\mathcal{J}_k^*(x_k) = \mathcal{J}_k(x_k)$ )

↳ obtained by DP algorithm

<pf within pf>  $\mathcal{J}_N^*(x_N) = g_N(x_N) = \mathcal{J}_N(x_N)$ .

Induction hypothesis: Assume for  $k+1$

$$\mathcal{J}_{k+1}^*(x_{k+1}) = \mathcal{J}_{k+1}(x_{k+1}) \quad \forall x_{k+1}.$$

$$\begin{aligned} \mathcal{J}_k^*(x_k) &= \min_{(\mu_k, \pi^{k+1})} E \left( g_N(x_N) + g_k(x_k, \mu_k(x_k), x_{k+1}) \right. \\ &\quad \left. + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1}) \right) \end{aligned}$$

$$\begin{aligned} &= \min_{\mu_k} E_{x_{k+1}} \left[ g_k(x_k, \mu_k(x_k), x_{k+1}) \right. \\ &\quad \left. + \min_{\pi^{k+1}} E_{x_{k+1} \mid x_k} \left( g_N(x_N) + \sum_{i=k+1}^{N-1} g_i(x_i, \mu_i(x_i), x_{i+1}) \right) \Big| x_{k+1} \right] \\ &\quad \text{↳ } \mathcal{J}_{k+1}^*(x_{k+1}) = \mathcal{J}_{k+1}(x_{k+1}) \text{ by induction hypothesis.} \end{aligned}$$

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

$$(*) \curvearrowright = \min_{a_k \in A(x_k)} E \left( g_k(x_k, \mu_k(x_k), x_{k+1}) + \mathcal{J}_{k+1}(x_{k+1}) \right) = \mathcal{J}_k(x_k)$$

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

(end of pf with pf)

To infer  $(\hat{x} \hat{x})$ , we need

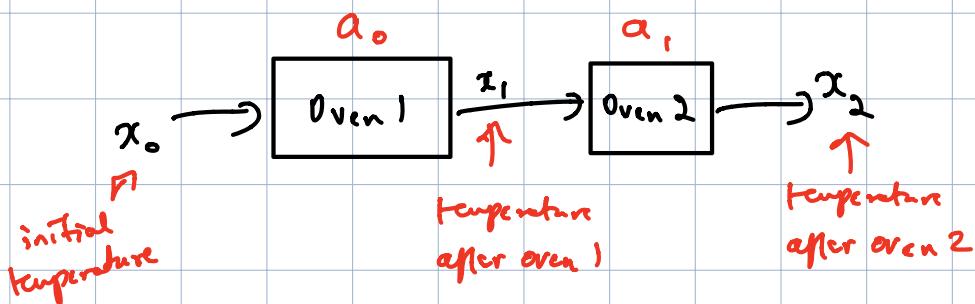
$$\min_{\mu \in B} F(x, \mu(x)) = \min_{a \in A(x)} F(x, a),$$

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

### Lecture - 3

#### Examples :

##### ① Linear System with quadratic cost



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

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

Evolution

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 :  $J_2(x_2) = (x_2 - T)^2$  — (xx)

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

from (x)  $\Rightarrow = \min_{a_1} (a_1^2 + J_2((1-\lambda)x_1 + \lambda a_1))$

wrong (xx)  $\Rightarrow = \min_{a_1} (a_1^2 + ((1-\lambda)x_1 + \lambda a_1 - T)^2)$

optimal action

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

$J_1(x_1) = \frac{((1-\lambda)x_1 - T)^2}{1+\lambda^2}$  ← quadratic

Check :

$\mu_0^*(x_0) = \frac{(1-\lambda)\lambda(T - (1-\lambda)^2 x_0)}{1+\lambda^2(1+(1-\lambda)^2)}$  ← linear

$J_0^*(x_0) = \frac{((1-\lambda)^2 x_0 - T)^2}{1+\lambda^2(1+(1-\lambda)^2)}$  ← quadratic

Adding randomness to ovens:

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

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

Applying DP algorithm:

$$\begin{aligned}
 J_1(x_1) &= \min_{a_1} E \left( a_1^2 + ((1-\lambda)x_1 + \lambda a_1 + \omega_1 - T)^2 \right) \\
 &= \min_{a_1} \left( a_1^2 + ((1-\lambda)x_1 + \lambda a_1 - T)^2 \right. \\
 &\quad \left. + 2E\omega_1((1-\lambda)x_1 + \lambda a_1 - T) \right. \\
 &\quad \left. + E\omega_1^2 \right) \\
 \stackrel{E\omega_1^2 \geq 0}{=} \min_{a_1} \left( a_1^2 + ((1-\lambda)x_1 + \lambda a_1 - T)^2 \right. &+ E\omega_1^2
 \end{aligned}$$

Minimizing RLS above leads to the same  
 $\mu_1^*(x_1)$  as before.

Lecture-4\*

Chess match - revisited (last time)

Consider an extension to N games  
 trivial  $\rightarrow P_d$ , bold  $\rightarrow P_w$

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

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

Apply DP algorithm:

$$J_k(x_k) = \max \left( P_d J_{k+1}(x_k) + (1-P_d) J_{k+1}(x_k - 1), \right)$$

(\*)  $\rightarrow$   
we are playing  
for rewards  
(choose min to max)  
 $\rightarrow$  DP algo

timid play

$$P_w J_{k+1}(x_k + 1) + (1-P_w) J_{k+1}(x_k - 1)$$

bold play

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

It is better to play bold when

$$\frac{P_w}{P_d} \rightarrow \frac{J_{k+1}(x_k) - J_{k+1}(x_k - 1)}{J_{k+1}(x_k + 1) - J_{k+1}(x_k - 1)}$$

$\rightarrow$  inferred  
from (\*)

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

$x_{N-1}$	$T_{N-1}$	Best play
$> 1$	1	does not matter
1	$\max(P_d + (1-P_d)P_w, P_w + (1-P_w)P_w)$ = $P_d + (1-P_d)P_w$	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  $T_{N-2}(0)$

$$\begin{aligned}
 T_{N-2}(0) &= \max \left( P_d P_w + (1-P_d) P_w^2, P_w (P_d + (1-P_d) P_w) \right. \\
 &\quad \left. + (1-P_w) P_w^2 \right) \\
 &= \max \left\{ P_w (P_d + (1-P_d) P_w), P_w (P_d + (1-P_d) P_w + (1-P_w) P_w) \right\} \\
 &= P_w (P_d + (1-P_d) P_w + (1-P_w) P_w) \Rightarrow \text{play bold}
 \end{aligned}$$

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

Another example: (Job scheduling)

N jobs to schedule

$T_i$  → 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" finished at time "t", then

The reward is  $\alpha^t R_i$ ,  $\alpha = \text{discount factor}$   
 $0 < \alpha < 1$

Cumulative reward = sum of each job's reward.

Goal: Schedule jobs to maximize cumulative reward.

"Interchange argument" to figure optimal schedule

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

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

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

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

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,

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

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  $\mu_{[1]} \geq \mu_{[2]} \geq \dots \geq \mu_{[N]}$

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

index-based

optimal policy

Further reading: Check out  
Gittin index  
Sec 1.3 of PLOC Vol.II

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's 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, finite mean

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

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

$A(s_0, x_0=0)$

$x_{k+1} = f(x_k, a_k, w_k)$ , 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  $E \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, a_k, w_k) \right]$   
 Expectation over  $w_0, w_1, \dots, w_{N-1}$

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

$$g_k(x_k, a_k, \omega_k) = \begin{cases} (1+r)^{N-k} x_k & \text{if } x_k \neq T, a_k = a' \\ 0 & \text{else} \end{cases}$$

not sold  
sell

$r > 0$  is the interest rate.

Apply DP algorithm:

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

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

$$\text{Let } \alpha_k = \frac{E(\tau_{k+1}(\omega_k))}{(1+r)^{N-k}}$$

Optimal policy: (threshold-based policy) "non-stationary thresholds  $\alpha_k$ "

sell if  $x_k > \alpha_k$

don't sell if  $x_k < \alpha_k$

if  $x_k = \alpha_k$   
both actions are fine

## Understanding the optimal policy:

Suppose

$$\alpha_k \geq \alpha_{k+1} \quad \forall k$$

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

DP algo in "V" notation is

$$V_N(x_N) = x_N$$

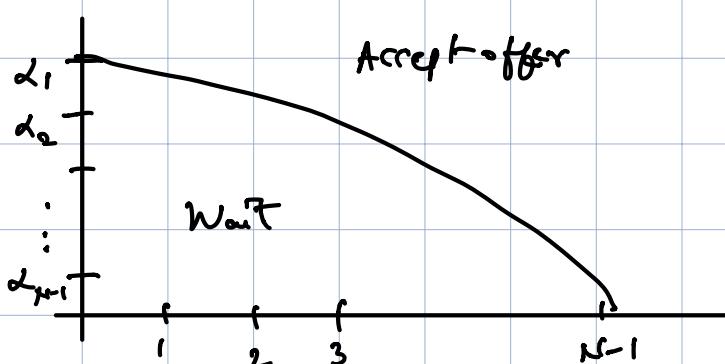
$$V_k(x_k) = \max \left( x_k, \frac{E(V_{k+1}(w_k))}{(1+r)} \right)$$

Optimal policy

$$\mu^*(x_k) = \begin{cases} a' & \text{if } x_k \geq \alpha_k \\ a'' & \text{else} \end{cases}$$

Lecture 5\*

Claim:  $\alpha_k \geq \alpha_{k+1}$



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

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 " $\omega$ " is a continuous, positive-valued r.v. with distribution  $F_\omega$  & density  $h$ .

$$V_{k+1}(\omega) = \begin{cases} \omega & \text{if } \alpha_{k+1} \leq \omega \\ \alpha_{k+1} & \text{else} \end{cases}$$

$$\alpha_k = \frac{E(V_{k+1}(\omega))}{1+\gamma}$$

$$= \frac{1}{1+\gamma} \int_0^\infty V_{k+1}(\omega) h(\omega) d\omega$$

$$= \frac{1}{1+\gamma} \left( \int_0^{\alpha_{k+1}} \alpha_{k+1} h(\omega) d\omega + \int_{\alpha_{k+1}}^\infty \omega h(\omega) d\omega \right)$$

(\*)  $\alpha_k = \frac{\alpha_{k+1}}{1+\gamma} F_\omega(\alpha_{k+1}) + \frac{1}{1+\gamma} \int_{\alpha_{k+1}}^\infty \omega h(\omega) d\omega$

$$\leq \frac{\alpha_{k+1} F_\omega(\alpha_{k+1})}{1+\gamma} + \frac{1}{1+\gamma} \int_0^\infty \omega h(\omega) d\omega$$

$$\alpha_k \leq \frac{\alpha_k}{1+\gamma} + \frac{E\omega}{1+\gamma}$$

From the above, we can conclude

$$0 \leq \alpha_k \leq \frac{E\omega}{\gamma}$$

Since  $\alpha_{k+1} \leq \alpha_k$ , the sequence  $\{\alpha_k\}$

converges as  $k \rightarrow \infty$ , say to some

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

Taking  $k \rightarrow -\infty$  in (x), we obtain

$$\bar{x} = \frac{\bar{x} F_w(\bar{x})}{1+\gamma}$$

$$+ \frac{1}{1+\gamma} \int_{\bar{x}}^{\infty} w h(w) dw$$

$\bar{x}$   
Can be calculated  
using  
distribution  
of  $w$

when  $N$  large

So, an approximation to the optimal policy is to  
use  $\bar{x}$  in place of  $x_k$  at stage  $k$ , i.e.,

Sell at  $x_k$  if  $x_k > \bar{x}$

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

} a threshold-based policy  
with "stationary" thresholds

<End of finite-horizon MDPs>