# Lec30

- Deterministic SSP: Principle of Optimality:

$$J_k(i) = \min_{j=1,2,\cdots,N} \left\{ a_{ij} + J_{k+1}(j) \right\}$$

- Stochastic DP

$$J_k(x_k) = \min_{u_k \in U_k(x_k)} E_{w_k} \left[ g_k(x_k, u_k, w_k) + J_{k+1}\left( f_k(x_k, u_k, w_k) \right) \right]$$

- Stopping time problems are another popular class of DP problems

- There is an action that needs to be taken

  - park a car

  - transmit a packet

- The decision is <u>when</u> to perform the action.

  - Once the action is performed, the system reaches the terminating state.

  - "Stopped".

- The problem is to figure out when to stop.

---

Ex 1)  Asset Selling

 - Bertsekas   P176

   Dynamic Programming and Optimal Control, Volume I, THIRD EDITION, Dimitri P. Bertsekas, Athena Scientific, Belmont, MA, 2005

- A person has an asset that needs to be sold before stage N.

- At each stage, he is offered $W_{k-1}$ amount of money for the asset.

  - future offers are random and independent.

  - Should he sell?

— Past offers expire immediately
— If he accepts the offer, he can invest the money at a fixed interest rate of $r$.
  — The payoff at the end (Stage $N$) is
$$W_{k-1}(1+r)^{N-k}$$
  — The last offer $W_{N-1}$ must be accepted if the asset has not been sold yet.

## Set up the DP

- States :
$$x_k = \begin{cases} T & , \text{sold ("terminating state")} \\ W_{k-1} & , \text{offer at stage } k \end{cases}$$

- Actions:
$$u_k = \begin{cases} 0 & , \text{do not sell} \\ 1 & , \text{sell} \end{cases}$$

- Transitions:
$$x_{k+1} = \begin{cases} T & , \text{if } x_k = T, \text{ or } x_k \neq T \text{ but } u_k = 1 \\ W_k & , \text{otherwise.} \end{cases}$$

- Payoff:
$$g_k = 0 \qquad \text{if} \qquad \text{do not sell}$$

$$g_k = \omega_{k-1}(1+r)^{N-k} \quad \text{if sell.}$$

$$g_N(x_N) = \begin{cases} x_N, & \text{if } x_N \neq T \\ 0, & \text{o/w.} \end{cases}$$

---

- Last stage:

  - Must sell

  - $J_N(x_N) = \begin{cases} x_N, & \text{if } x_N \neq T \\ 0, & \text{o/w.} \end{cases}$

- For $k \leq N$ stage:

  $$J_k(x_k) = \begin{cases} \max\left[ (1+r)^{N-k} x_k, \; E\left[ J_{k+1}(\omega_k) \right] \right] \overset{\downarrow \omega_{k-1}}{} \\ \qquad\qquad\qquad\qquad \text{if } x_k \neq T \\ 0 \qquad\qquad\qquad\qquad \text{if } x_k = T \end{cases}$$

- Thus, the optimal policy is to accept an offer if it is greater than

  $$\frac{E\left[ J_{k+1}(\omega_k) \right]}{(1+r)^{N-k}} \overset{\circ}{=} \alpha_k$$

  which can be viewed as the expected revenue of future offers discounted to the present time.

- Note that this is an example of a threshold policy

  - Take an action if the input is larger than a threshold.

- In this problem, the existence of an optimal threshold policy is quite intuitive

  - If I decide to accept an offer at $w_{k-1} = a$, I would also accept offers $\geq a$.

  - True if future offers are independent of past offers.

## Properties of the optimal thresholds

- Assume that the offers $w_k$ are i.i.d.

- We will show that the thresholds are decreasing

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

  - In other words, the owner is more willing to sell as the deadline comes closer.

  - Alternatively, if an offer is good enough at $k$, it will also be good enough at a later time when there are fewer chances of good offers.

- Recall that

$$\alpha_k = \frac{E[J_{k+1}(w)]}{(1+r)^{N-k}} \quad \overset{i.i.d}{\checkmark}$$

then

$$\alpha_{k-1} = \frac{E[J_k(w)]}{(1+r)^{N-k+1}}$$

To show $\alpha_{k-1} \geq \alpha_k$, it is sufficient to show that

$$\frac{J_k(x)}{(1+r)^{N-k+1}} \geq \frac{J_{k+1}(x)}{(1+r)^{N-k}}$$

for all $x$.

$$\Longrightarrow \quad \frac{J_k(x)}{1+r} \geq J_{k+1}(x)$$

— For $k = N-1$,

$$J_{k+1}(x) = J_N(x) = x$$

$$J_k(x) = J_{N-1}(x) = \max \left\{ (1+r)x, \; \bar{E}[J_N(\omega)] \right\}$$

$$\geq (1+r)x = J_{k+1}(x) \cdot (1+r)$$
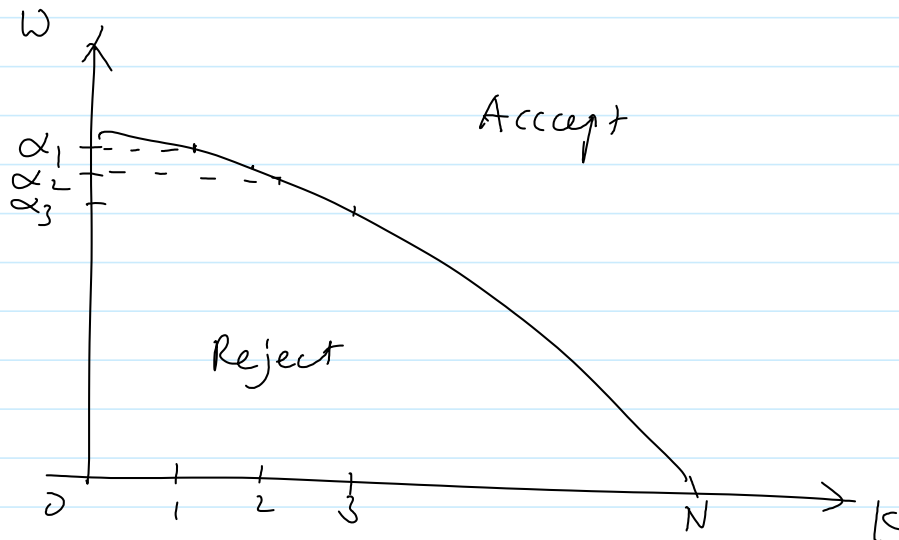
— Assume that

$$\frac{J_k(x)}{1+r} \geq J_{k+1}(x) \qquad \text{for all } x$$

then

$$\frac{J_{k-1}(x)}{1+r}$$

$$= \frac{1}{1+r} \max \left\{ x(1+r)^{N-k+1}, \; \bar{E}[J_k(\omega)] \right\}$$

$$\geq \frac{1}{1+r} \max \left\{ x(1+r)^{N-k+1}, \; \bar{E}[(1+r)J_{k+1}(\omega)] \right\}$$

$$= \max \left\{ x(1+r)^{N-k}, \; \bar{E}[J_{k+1}(\omega)] \right\}$$

$$= J_k(x)$$

— Done!

W ↑

Accept

$\alpha_1$ - - - - - - -
$\alpha_2$ - - - - - - -
$\alpha_3$ - - - - - - -

Reject

0    1    2    3                  N → k

- Can also show that $\alpha_k \to \bar{\alpha}$ as $k \to +\infty$

  - A stationary policy is optimal for infinite horizon

- May also be generalized to the case where the offers are temporally correlated and are generated by a 1-st order AR process.

# Complexity

- Suppose that each minimization costs $A$

- At each stage $N, N-1, ..$

  - There are $S$ possible states

- Total cost: $N \cdot S \cdot A$

---

- Compare this with the earlier $LP$ formulation

- The # of variables $y_{sa}^k$ is already $N \cdot S \cdot A$

- Thus, exploiting the DP structure allows us to build a much efficient solution!

---

"Curse - of - dimensionality"

- Linear in the state space

- If each state can be represented by $L$ dimension, each dimension has $b$ values

  - Total state space is $b^L$

  - Quickly become prohibitive as $L$ increases

- We will discuss possible approach at the end.

- Let us see another example of optimal stopping in wireless networks

    - N. B. Chang, M. Liu, "Optimal Channel Probing and Transmission Scheduling for Opportunistic Spectrum Access," IEEE/ACM Transactions on Networking, Dec. 2009

- A wireless system has $N$ channels

- The transmitter wishes to pick one channel for transmission.

- Each channel $j$ has quality (payoff) $X_j$

    - random with known distribution

    - independent of other channels

- However, in order to know the value of $X_j$, the xmitter must probe the channel $j$, which incurs $C_j > 0$

## Optimal Stopping Problem

- $N$ stages

- At stage $j$, let $S$ denote the set of unprobed channels

    - Channels in $\Omega - S$ have been probed

— Their quality $X_j$, $j \in \Omega - S$, is known

— The xmitter needs to decide

① probe a new channel in $S$

stop
{
② use one previously probed channel to transmit : "retire"

③ use a channel in $S$ to transmit : "guess"
}

— Assume channel qualities are fixed until a transmission takes place.

## Goal :

— To maximize net payoff

$$J^* = \max \ \overline{E}\left[ X_{\pi(\tau)} - \sum_{t=1}^{\tau-1} C_{\pi(t)} \right]$$

where $z(1), z(2), \cdots, z(\tau)$ denotes the sequence of channels probed until transmission

## Optimality Equation

— State at a particular stage :

— $\Omega - S$ : set of channels probed

— The quality of these channels. However, we only need to remember the best channel

$$u = \max \{ X_j \mid j \in \Omega - S \}$$

— State is $(u, S)$

— Let $J(u, S)$ be the max cost to go from state $(u, S)$

- If $S = \phi$ : $J(u, S) = u$

Only action is "retire"

— for a given S, Depending on the action:

① probe $j \in S$ : payoff $-c_j$

next state : $S - j$

$$\max \{ u, X_j \}$$

② retire : payoff $u$

③ guess $j \in S$ : payoff $E[X_j]$

$$J(u, S) = \max \left\{ \max_{j \in S} \left[ -c_j + E\left[ J\left( \max \{u, X_j\}, S-j \right) \right] \right] \right.$$

$$u,$$

$$\left. \max_{j \in S} E[X_j] \right\}$$

— Start from $(u, S = \Omega)$ :

$$J(n, \Omega) = u$$

— Do backward induction. State space: exp in $N$

---

# Challenge

- The value of $u$ is continuous!

- Need discretization as an approximation.

- Deeper understanding of the structure of the optimal policy will be helpful

# Threshold Property

- $J(n, s)$ must be non-decreasing in $u$.

- If $J(u, s) = u$, then for any $\hat{n} \geq u$, we must have $J(\hat{n}, s) = \hat{n}$

- If $J(u, s) = E(X_j)$, then for any $\hat{n} \leq u$, we must have $J(\hat{n}, s) = E(X_j)$.

- These properties imply that, for fixed $s$, the optimal policy has a threshold structure w.r.t. $u$

  In particular, define
  - $a_s = \inf \{u : J(u, s) = u\}$

  - $b_s = \sup \{u : J(u, s) = E[X_j], \text{ for some } j \in s\}$

  Then: $0 \leq b_s \leq a_s \leq M$. The optimal

policy must be

$$z^*(n, s) = \begin{cases} \text{retire}(n) & \text{if } u \geq a_s \\ \text{probe}(jn), jn \in S & \text{if } b_s < u < a_s \\ \text{guess}(j), j \in S & \text{if } u < b_s \end{cases}$$

— The idea is then to iterate over these thresholds

Start from the last stage:

— $S = \{j\}$ : only one channel remains to be probed.

— $J(u, \{j\}) =$

$$\max \begin{cases} -c_j + E[\max(u, X_j)], \\ u, \\ E(X_j) \end{cases} \quad \begin{matrix} \text{probe} \\ \text{retire} \\ \text{guess} \end{matrix}$$

— To determine $a_j$

$$a_j = \min \{ u: \quad u \geq E[X_j], $$

$$\underbrace{u \geq -c_j + E[\max(u, X_j)]}_{\Updownarrow} \}$$

$$c_j \geq E[(X_j - u)^+]$$

— To determine $b_j$

$$b_j = \max \{ u: \quad u \leq E(X_j),$$

$$\underbrace{E(X_j) \geq -c_j + E[\max(u, X_j)]}_{\Updownarrow} \}$$

$$c_j \geq E[(u - X_j)^+]$$

expected reward

retire
probe : case 1
probe : case 2

$\overline{E}(X_j)$

$a_j$

$b_j$

guess

$E[\max(X_j, w)] - C_j$

$b_j = a_j = E(X_j)$

$u$

— Note: the two cases of "probe" are parallel.

— Thm: $C_j$ determines the gap between $a_j$ & $b_j$.

— Other steps of the iteration can be determined in a similar manner, although more involved.

- Linear - Quadratic Problems

$$x_{k+1} = A_k x_k + B_k U_k + W_k$$

- Inventory Control

$$x_{k+1} = x_k + U_k - W_k$$

         ↑    ↑

     purchase  demand

- Scheduling and interchange arguments

---

<u>Problems with</u> deterministic policies as solution

- Suppose that we have N jobs.

- Job $i$ needs $T_i$ amount of time to complete

    - $T_i$'s are random but independent.

- If Job $i$ is completed at time $t$, the reward is $\alpha^t R_i$, with $0 < \alpha < 1$.

    - $R_i$ are given constants

    - Wish to complete more valuable jobs first.

- The problem is to find a schedule that maximizes the total expected reward.

---

- It is easy to see that, given a subset of jobs

yet to be completed, the future schedule is independent from the past.

- True since $T_i$'s are independent, and all rewards are scaled by $\alpha^t$.

— It is then clear that the optimal policy can be mapped to a deterministic schedule.
$$(i_0, i_1, \dots, i_{N-1})$$

— In order to find the optimal deterministic schedule, we may use the following interchange argument.

- Suppose that $i, j$ are two subsequent jobs in a schedule.

- Let us see if interchanging them will be beneficial or not.

- Let $t_0$ be the completion time for the previous job before $i$ & $j$

$$(i, j) \rightarrow E\left[\alpha^{t_0 + T_i} R_i + \alpha^{t_0 + T_i + T_j} R_j\right] \overset{\Delta}{=} A$$

$$(j, i) \rightarrow E\left[\alpha^{t_0 + T_j} R_j + \alpha^{t_0 + T_i + T_j} R_i\right] \overset{\Delta}{=} B$$

- Since $T_i, T_j$, and $t_0$ are independent

$$A \geq B$$

$$\Leftrightarrow \quad R_i \, E\left[\alpha^{T_i}\right] + R_j \, E\left(\alpha^{T_i}\right) E\left(\alpha^{T_j}\right)$$

$$\geq R_j \, E\left[\alpha^{T_j}\right] + R_i \, E\left[\alpha^{T_i}\right] E\left(\alpha^{T_j}\right)$$

$$\Leftrightarrow \quad R_i \, E\left(\alpha^{T_i}\right)\left[1 - E\left[\alpha^{T_j}\right]\right]$$

$$\geq R_j \, E\left(\alpha^{T_j}\right)\left[1 - E\left(\alpha^{T_i}\right)\right]$$

$$(\Rightarrow) \qquad \frac{R_i \, \mathcal{E}(\alpha^{T_i})}{1 - \mathcal{E}(\alpha^{T_i})} \geq \frac{R_j \, \mathcal{E}(\alpha^{T_j})}{1 - \mathcal{E}(\alpha^{T_j})}$$

— The optimal deterministic schedule should be
decreasing in

$$\frac{R_i \, \mathcal{E}(\alpha^{T_i})}{1 - \mathcal{E}(\alpha^{T_i})}$$

— In general, such an interchanging argument
may not always work. Nonetheless, it can be
the basis for useful heuristics