

# Lec34

Sunday, April 16, 2023 5:27 PM

Are you okay with posting your report on a password-protected website?

- Please complete course evaluation:
- Course feedback form
- Final project presentation:
  - o Room?
  - o Watch out for email announcements

- Deterministic SSP: Principle of Optimality:

$$J_K(i) = \min_{j=1,2,\dots,N} \{ a_{ij} + J_{K+1}(j) \}$$

- Finite horizon stochastic SP

$$J_K(x_K) = \min_{u_K \in U_K(x_K)} \mathbb{E}_{w_K} \left[ g_K(x_K, u_K, w_K) + J_{K+1}(f_K(x_K, u_K, w_K)) \right]$$

- Infinite-horizon SSP

$$J^*(i) = \min_u \left\{ g(i, u) + \sum_j P_{ij}(u) J^*(j) \right\}$$

- Discounted problems

- Using the SSP mapping

$$J^*(i) = \min_u \left\{ g(i, u) + \sum_j \underbrace{\alpha P_{ij}(u)}_{\substack{\text{transition} \\ \text{probability} \\ \text{in SSP}}} J^*(j) \right\}$$

$$\Leftrightarrow J^*(i) = \min_u \left\{ g(i, u) + \alpha \sum_j P_{ij}(u) J^*(j) \right\}$$

↑  
future cost  
is discounted  
by  $\alpha$

↑  
future cost  
from  $j$

- Average-cost problem

$$J_{\lambda}(i) = \lim_{N \rightarrow +\infty} \frac{1}{N} \mathbb{E} \left\{ \sum_{k=0}^{N-1} g(x_k, \mu_k(x_k)) \mid x_0 = i \right\}$$

$$\lambda^* = \min_{\lambda} J_{\lambda}(i)$$

$$\lambda^* + h(i) = \min_u \left\{ g(i, u) + \sum_j p_{ij}(u) h(j) \right\}$$

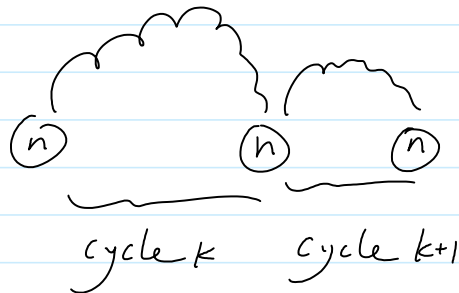
## Connection to SSP

Sunday, March 29, 2015 10:03 AM

Assumption:

- finite state space
  - bounded cost per-stage
  - "Exponential recurrent": There is one state  $n$  such that for some integer  $m > 0$ , and for all initial states and all policies, the state  $n$  is visited with positive probability at least once within the first  $m$  stages
    - let  $p$  be the minimum probability of not entering state  $n$  at least once in  $m$  stages
    - The probability of not entering state  $n$  in  $mk$  stages goes down as  $p^k$
- 

- We can then think of the infinite-horizon problem into cycles of successive visit to the state  $n$ .



- All such cycles are statistically the same (i.i.d).
  - start/end at the same state
  - same transition probabilities

- Intuitively, if we optimize an appropriate "average cost" in each of these cycles, we will be able to optimize the average cost for the entire horizon

What should we optimize in each cycle?

- Let  $C_{nn}(\mu) = \text{cost from } n \text{ to } n$   
 $N_{nn}(\mu) = \text{time from } n \text{ to } n$

- Should we optimize  $E\left[\frac{C_{nn}}{N_{nn}}\right] \stackrel{?}{=} \lambda'$ ?

The answer is no, because  $\lambda'$  differs from the average cost for the entire horizon (i.e., averaged over all cycles)?

$$\frac{C_{nn}^1 + C_{nn}^2 + \dots + C_{nn}^k}{N_{nn}^1 + N_{nn}^2 + \dots + N_{nn}^k}$$

$$= \frac{\frac{C_{nn}^1 + C_{nn}^2 + \dots + C_{nn}^k}{k}}{\frac{N_{nn}^1 + N_{nn}^2 + \dots + N_{nn}^k}{k}}$$

$$\rightarrow \text{as } k \uparrow \quad \frac{E[C_{nn}]}{E[N_{nn}]} \neq E\left[\frac{C_{nn}}{N_{nn}}\right]$$

↑  
we should optimize this instead!

- However,  $\frac{E[C_{nn}]}{E[N_{nn}]}$  is not an additive cost!

- Let  $\lambda^*$  = optimal average cost per-stage.

- Note that

$$\frac{E[C_{nn}(\mu)]}{E[N_{nn}(\mu)]} \geq \lambda^* \quad \text{for all policy } \mu$$

$$\Rightarrow E[C_{nn}(\mu) - N_{nn}(\mu) \cdot \lambda^*] \geq 0$$

with equality attained if  $\mu$  is optimal.

- Now consider an SSP with per-stage cost

$$g(i, n) - \lambda^*$$

and terminating at  $n$ .

- The total cost is exactly

$$E[C_{nn}(\mu) - N_{nn}(\mu) \lambda^*]$$

- Any policy will produce such a total cost  $\geq 0$

- But the optimal policy will make it 0!

$\Rightarrow$  The optimal policy  $\mu$  will also be the optimal for the SSP

- with  $J^*(n) = 0$

---

## Bellman's Equation

- We now use the SSP to derive Bellman's equation for the average cost problem.

- Let  $h^*(i)$  be the optimal cost-to-go for this SSP problem, starting from state  $i$

$$h^*(n) = 0$$

- Bellman's Equation

$$h^*(i) = \min_u \left[ g(i, u) - \lambda^* + \sum_j P_{ij}(u) h^*(j) \right]$$

or

$$h^*(i) + \lambda^* = \min_u \left[ g(i, u) + \sum_j P_{ij}(u) h^*(j) \right] \quad (*)$$

- All results follow from that of SSP
- Starting from any  $(h_0(i))$ , the DP iteration

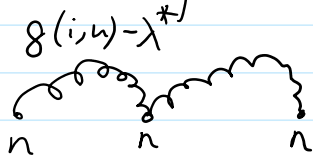
$$h^{k+1}(i) + \lambda^* = \min_u \left[ g(i, u) + \sum_j P_{ij}(u) h^k(j) \right]$$

will converge to  $h^*(j)$

- $h^*(j)$  satisfies the above equation (\*) and is unique.
- any policy that minimizes the RHS is optimal.

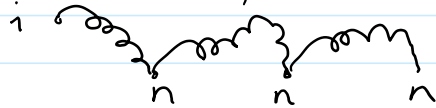
What is the meaning of  $h^*(i)$ ?

- If we start from  $n$ , then  $h^*(n) = 0$



- The cost accumulated exactly cancels out with  $\lambda^*$
- If we start from  $i$ ,

- If we start from  $i$ ,



$h^*(i)$  captures the difference between the cost accumulated and  $\lambda^*$   
 $\Rightarrow$  "relative" value function.

---

- The problem, however, is we do not know  $\lambda^*$  yet!

- Fortunately, since  $h^*(n) = 0$ , there are exactly  $(n-1)$  unknown  $h^*(i)$  + 1 unknown  $\lambda^*$

- Further, since we can also write an equation for  $i=n$ , we have a total of  $n$  equations.

- Hence, the solution to (\*) is likely unique (by restricting  $h^*(n) = 0$ )

- Is it always the case?

## Example

Saturday, April 25, 2015 2:07 PM

- Bertsekas P429
  - A "lazy" worker receives a new order in each period with probability  $p$ , independently of other periods
  - However, she does not want to work whenever a new order arrives, because she is very efficient at batch-processing
  - Rather, if she waits and processes all orders in a batch, she only incurs one set-up cost of  $K > 0$ .
  - On the other hand, the cost for each unfilled order at each period is  $c > 0$
  - Assume that the max # of unfilled order is  $n$ , in which case the worker must process them
  - What is the policy that minimizes the average cost?
- 

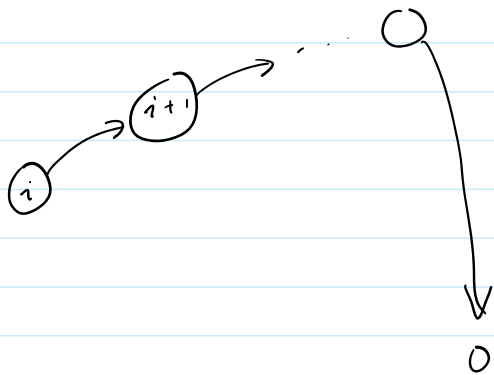
### Bellman's Equation

- state: # of unfilled orders  $i$
- actions:
  - process: cost  $k$ 
    - next state: 0 or 1
  - wait: cost  $c_i$



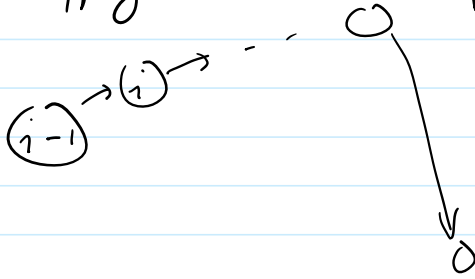
- next state  $i$  or  $i+1$
- must process if  $i = n$ .
- Bellman's Equation
 
$$h^*(i) + \lambda = \min \left\{ \begin{array}{l} K + p h^*(1) + (1-p) h^*(0), \\ c_i + p h^*(i+1) + (1-p) h^*(i) \end{array} \right\}_{i=0, 1, \dots, n-1}$$

$$h^*(n) + \lambda = K + p h^*(1) + (1-p) h^*(0)$$
- The policy is to process the orders if
 
$$c_i + p h^*(i+1) + (1-p) h^*(i) \geq K + p h^*(1) + (1-p) h^*(0)$$
- It is intuitive that  $h^*(i)$  is non-decreasing in  $i$ 
  - $h^*(i)$  is the cost to go from state  $i$  for the SSP problem to a recurrent state, say, 0.
  - The cost each step is  $c_i$  or  $K$ , which is non-decreasing in  $i$ .



- Suppose that a sequence of decisions is optimal, which will attain  $h^*(i)$
- Apply the same sequence of decisions to

- which will attain  $h^*(i)$
- Apply the same sequence of decisions to



- The expected cost should only be lower
- Since  $h^*(i-1)$  corresponds to the minimum cost for the SSP, it will be even lower

$$\Rightarrow h^*(i-1) \leq h^*(i)$$

- Thus, the optimal policy must have a threshold structure:

There exists  $i_0$  such that if  $i \geq i_0$ , process all unfilled orders.

- It remains to answer the question what the resulting  $\lambda^*$  for the Bellman's Equation is always correct.
  - Another question is what is the meaning of  $h^*(i)$ ?
- 

Proposition: (Bertsekas P426)

- (a) If a scalar  $\lambda$  and a vector  $h = [h(1) \dots h(n)]$  satisfies Bellman's Equation

$$\lambda + h(i) = \min_u \left\{ f(i, u) + \sum_j P_{ij}(u) h(j) \right\} \quad (*)$$

Then  $\lambda$  is the optimal average cost starting from any state  $i$ , and such  $h(j)$  is unique given  $h(n) = 0$

- (b) For any stationary policy  $\mu$ , there exists a unique vector  $[h_\mu(1) \dots h_\mu(n)]$  with  $h_\mu(n) = 0$ , and a unique  $\lambda_\mu$  such that

$$\lambda_\mu + h_\mu(i) = f(i, \mu(i)) + \sum_j P_{ij}(\mu(i)) h_\mu(j)$$

- (c) A stationary policy is optimal if & only if it maximizes the RHS of the Bellman's Equation.
- 

Proof of (a):

Suppose

$$\lambda + h(i) = \min_u \left\{ f(i, u) + \sum_j P_{ij}(u) h(j) \right\}$$

Consider a  $k$ -stage problem with the terminating cost being  $h(i)$ .

- The minimum one-stage cost is

$$J_1(i) = \min_n \left\{ f(i, n) + \sum_j P_{ij}(n) h(j) \right\}$$
$$= \lambda + h(i) \quad \text{for all } i$$

- The minimum 2-stage cost is

$$J_2(i) = \min_n \left\{ f(i, n) + \sum_j P_{ij}(n) J_1(j) \right\}$$
$$= \lambda + \min_n \left\{ f(i, n) + \sum_j P_{ij}(n) h(j) \right\}$$
$$= 2\lambda + h(i)$$

- By induction, we can show that the  $\min$   $k$ -stage cost is  $k\lambda + h(i)$

- As  $k \rightarrow +\infty$ , the min average cost  $\rightarrow$  go must be  $\lambda$

- Finally,  $h(i)$  is unique because (\*) is identical to the Bellman's equation for the SSP problem to state  $n$ .

- The solution to the latter problem is unique.

### Proof of (b)

- Think of  $h_{\mu}(i)$  as the cost  $\rightarrow$  go of the SSP problem with termination state  $n$   
\* per-stage cost of

$$g(i, \mu(i)) - \lambda_{\mu}$$

Then, this equation in part (b) must hold.

$\lambda_{\mu}$  is unique because  $\lambda_{\mu}$  must be the average cost of the policy, which can be shown as in part (a)

$h_{\mu}(i)$  is unique since there are  $n$  equations &  $n$  variables.

Proof of (c)

Compare the two equations.

## Value iteration

Saturday, April 25, 2015 3:27 PM

- To use Bellman's equation, one can solve it directly
  - May be involved
- Or, use the following "natural" value iteration.
  - Start from any initial  $J_0(i) \dots J_0(n)$
  - Use the DP iteration to get the min  $(k+1)$ -stage cost

$$J_{k+1}(i) = \min_u g(i,u) + \sum_j P_{ij}(u) J_k(j)$$

- Based on our analysis, if  $J_0(i) = h^*(i)$ , then

$$J_{k+1}^*(i) = (k+1)\lambda^* + h^*(i)$$

For other values of  $J_0(i)$ , the difference between  $J_{k+1}(i)$  &  $J_{k+1}^*(i)$  is at most

$$\max_i |J_0(i) - h^*(i)|$$

because the only difference is the terminal cost.

- Hence, regardless of  $J_0(i)$ , we must have

$$\frac{J_k(i)}{k} \rightarrow \lambda^*$$

- However, numerically this method runs into difficulties where  $J_k(i) \rightarrow +\infty$  as  $k \rightarrow +\infty$
- Further, it does not tell us the value of  $h^*(i)$ .
  - Why is it not  $J_k(i) - \lambda^* k$ ?

## Relative Value Iteration

- We can subtract any <sup>fixed</sup> value from  $J_k(i)$  for all  $i$ .  
It does not change the min operation in the DP iteration.
- One way is to subtract a value so that one element, denoted by  $h_k(n)$  is always zero.

$$h_{k+1}(i) = \min_n \left\{ g(i, n) + \sum_j P_{ij}(n) h_k(j) \right\}$$

$$= \min_n \left\{ g(n, n) + \sum_j P_{nj}(n) h_k(j) \right\}$$

- Then, we can show that  $h_k(i) \rightarrow h^*(i)$   
(with an additional assumption)

## Policy iteration

Saturday, April 25, 2015 3:53 PM

- Alternately, we can use policy iteration
- Start with any stationary policy  $\mu^0$
- For policy  $\mu^k$ , find the average cost by solving

$$h^k(i) + \lambda^k = g(i, \mu^k(i)) + \sum_j P_{ij}(\mu^k(i)) h^k(j)$$

- with  $h^k(n) = 0$

- a linear program.

- Policy improvement

$$\mu^{k+1}(i) = \operatorname{argmin}_\mu g(i, \mu) + \sum_j P_{ij}(\mu) h^k(j)$$

- We can show that (Bertsekas p433)

either  $\lambda^{k+1} < \lambda^k$

or  $\lambda^{k+1} = \lambda^k$  &  $h^{k+1}(i) \leq h^k(i) \quad \forall i$

- Since there are a finite # of stationary policies, this method must converge to the optimal policy in a finite # of steps.



The average-cost problem also has a connection to a linear program!

Describe a stationary policy as follows

- Let  $\lambda_{iu}$  be the steady-state prob. of being at state  $i$  and taking action  $u$ .

$$\sum_i \sum_u \lambda_{iu} = 1$$

- Since it is the steady-state prob., it should also satisfy a balance equation

$$\sum_u \lambda_{iu} = \sum_j \sum_u \lambda_{ju} \cdot P_{ji}(u) \quad \text{for all } i.$$

- The average-cost is given by

$$\sum_{iu} \lambda_{iu} \cdot g(i, u)$$

- Hence, the average cost problem can be rewritten as the following linear program

$$\begin{aligned} \min \quad & \sum_{iu} \lambda_{iu} g(i, u) \\ \text{sub to} \quad & \sum_u \lambda_{ju} = \sum_{iu} \lambda_{iu} P_{ij}(u) \quad \forall j \\ & \sum_{iu} \lambda_{iu} = 1 \end{aligned}$$

- Our analysis earlier shows that a deterministic policy is optimal, i.e. for each state  $i$ , only one  $u$  needs to have

$$\lambda_{iu} \neq 0$$

- Not true in constrained MDP.

- Please complete course evaluation:
- Course feedback form
- Final project presentation:
  - o Room?
  - o Watch out for email announcements