

6.231 DYNAMIC PROGRAMMING

LECTURE 7

LECTURE OUTLINE

- DP for imperfect state info
- Sufficient statistics
- Conditional state distribution as a sufficient statistic
- Finite-state systems
- Examples

REVIEW: IMPERFECT STATE INFO PROBLEM

- Instead of knowing x_k , we receive observations

$$z_0 = h_0(x_0, v_0), \quad z_k = h_k(x_k, u_{k-1}, v_k), \quad k \geq 0$$

- I_k : information vector available at time k :

$$I_0 = z_0, \quad I_k = (z_0, z_1, \dots, z_k, u_0, u_1, \dots, u_{k-1}), \quad k \geq 1$$

- Optimization over policies $\pi = \{\mu_0, \mu_1, \dots, \mu_{N-1}\}$, where $\mu_k(I_k) \in U_k$, for all I_k and k .
- Find a policy π that minimizes

$$J_\pi = \underset{\substack{x_0, w_k, v_k \\ k=0, \dots, N-1}}{E} \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(I_k), w_k) \right\}$$

subject to the equations

$$x_{k+1} = f_k(x_k, \mu_k(I_k), w_k), \quad k \geq 0,$$

$$z_0 = h_0(x_0, v_0), \quad z_k = h_k(x_k, \mu_{k-1}(I_{k-1}), v_k), \quad k \geq 1$$

DP ALGORITHM

- DP algorithm:

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

for $k = 0, 1, \dots, N - 2$, and for $k = N - 1$,

$$J_{N-1}(I_{N-1}) = \min_{u_{N-1} \in U_{N-1}} \left[E_{x_{N-1}, w_{N-1}} \left\{ g_{N-1}(x_{N-1}, u_{N-1}, w_{N-1}) \right. \right. \\ \left. \left. + g_N(f_{N-1}(x_{N-1}, u_{N-1}, w_{N-1})) \mid I_{N-1}, u_{N-1} \right\} \right]$$

- The optimal cost J^* is given by

$$J^* = E_{z_0} \{ J_0(z_0) \}.$$

SUFFICIENT STATISTICS

- Suppose there is a function $S_k(I_k)$ such that the min in the right-hand side of the DP algorithm can be written in terms of some function H_k as

$$\min_{u_k \in U_k} H_k(S_k(I_k), u_k)$$

- Such a function S_k is called a **sufficient statistic**.
- An optimal policy obtained by the preceding minimization can be written as

$$\mu_k^*(I_k) = \bar{\mu}_k(S_k(I_k)),$$

where $\bar{\mu}_k$ is an appropriate function.

- Example of a sufficient statistic: $S_k(I_k) = I_k$
- Another important sufficient statistic

$$S_k(I_k) = P_{x_k|I_k},$$

assuming that v_k is characterized by a probability distribution $P_{v_k}(\cdot | x_{k-1}, u_{k-1}, w_{k-1})$

DP ALGORITHM IN TERMS OF $P_{x_k|I_k}$

- **Filtering Equation:** $P_{x_k|I_k}$ is generated recursively by a dynamic system (estimator) of the form

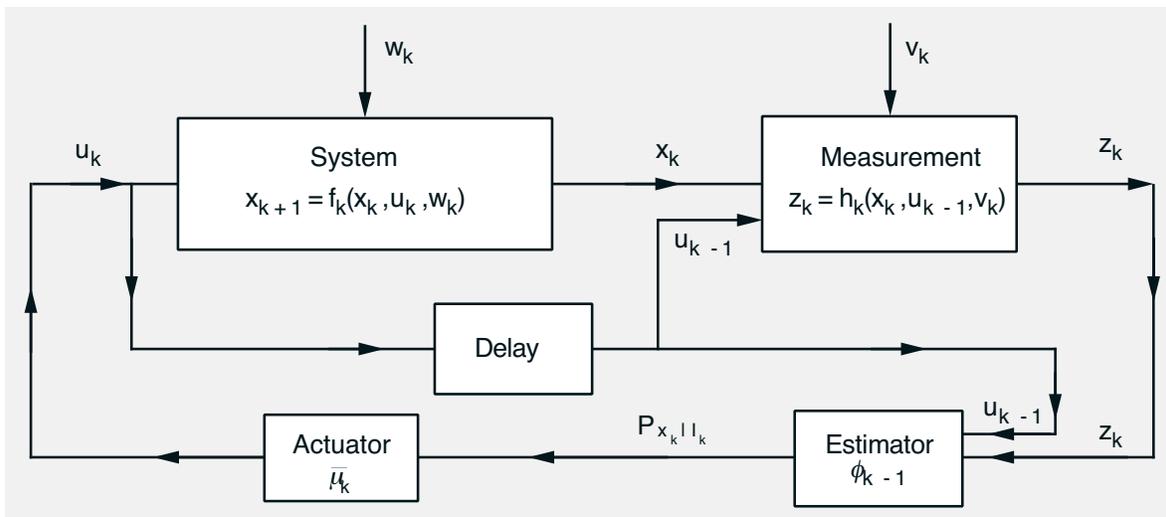
$$P_{x_{k+1}|I_{k+1}} = \Phi_k(P_{x_k|I_k}, u_k, z_{k+1})$$

for a suitable function Φ_k

- DP algorithm can be written as

$$\bar{J}_k(P_{x_k|I_k}) = \min_{u_k \in U_k} \left[E_{x_k, w_k, z_{k+1}} \left\{ g_k(x_k, u_k, w_k) + \bar{J}_{k+1}(\Phi_k(P_{x_k|I_k}, u_k, z_{k+1})) \mid I_k, u_k \right\} \right]$$

- It is the DP algorithm for a **new problem** whose state is $P_{x_k|I_k}$ (also called **belief state**)



EXAMPLE: A SEARCH PROBLEM

- At each period, decide to search or not search a site that may contain a treasure.
- If we search and a treasure is present, we find it with prob. β and remove it from the site.
- Treasure's worth: V . Cost of search: C
- States: treasure present & treasure not present
- Each search can be viewed as an observation of the state
- Denote

p_k : prob. of treasure present at the start of time k
with p_0 given.

- p_k evolves at time k according to the equation

$$p_{k+1} = \begin{cases} p_k & \text{if not search,} \\ 0 & \text{if search and find treasure,} \\ \frac{p_k(1-\beta)}{p_k(1-\beta)+1-p_k} & \text{if search and no treasure.} \end{cases}$$

This is the **filtering equation**.

SEARCH PROBLEM (CONTINUED)

- DP algorithm

$$\bar{J}_k(p_k) = \max \left[0, -C + p_k \beta V \right. \\ \left. + (1 - p_k \beta) \bar{J}_{k+1} \left(\frac{p_k(1 - \beta)}{p_k(1 - \beta) + 1 - p_k} \right) \right],$$

with $\bar{J}_N(p_N) = 0$.

- Can be shown by induction that the functions \bar{J}_k satisfy

$$\bar{J}_k(p_k) \begin{cases} = 0 & \text{if } p_k \leq \frac{C}{\beta V}, \\ > 0 & \text{if } p_k > \frac{C}{\beta V}. \end{cases}$$

- Furthermore, it is optimal to search at period k if and only if

$$p_k \beta V \geq C$$

(expected reward from the next search \geq the cost of the search - a **myopic rule**)

FINITE-STATE SYSTEMS - POMDP

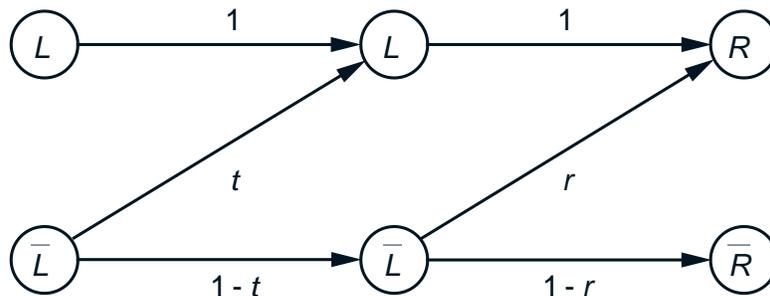
- Suppose the system is a finite-state Markov chain, with states $1, \dots, n$.
- Then the conditional probability distribution $P_{x_k|I_k}$ is an n -vector

$$\left(P(x_k = 1 | I_k), \dots, P(x_k = n | I_k) \right)$$

- The DP algorithm can be executed over the n -dimensional simplex (state space is not expanding with increasing k)
- When the control and observation spaces are also finite sets the problem is called a POMDP (Partially Observed Markov Decision Problem).
- For POMDP it turns out that the cost-to-go functions \bar{J}_k in the DP algorithm are piecewise linear and concave (Exercise 5.7)
- Useful in practice both for exact and approximate computation.

INSTRUCTION EXAMPLE I

- Teaching a student some item. Possible states are L : Item learned, or \bar{L} : Item not learned.
- **Possible decisions:** T : Terminate the instruction, or \bar{T} : Continue the instruction for one period and then conduct a test that indicates whether the student has learned the item.
- **Possible test outcomes:** R : Student gives a correct answer, or \bar{R} : Student gives an incorrect answer.
- Probabilistic structure



- **Cost of instruction:** I per period
- **Cost of terminating instruction:** 0 if student has learned the item, and $C > 0$ if not.

INSTRUCTION EXAMPLE II

- Let p_k : prob. student has learned the item given the test results so far

$$p_k = P(x_k = L \mid z_0, z_1, \dots, z_k).$$

- **Filtering equation:** Using Bayes' rule

$$\begin{aligned} p_{k+1} &= \Phi(p_k, z_{k+1}) \\ &= \begin{cases} \frac{1-(1-t)(1-p_k)}{1-(1-t)(1-r)(1-p_k)} & \text{if } z_{k+1} = R, \\ 0 & \text{if } z_{k+1} = \bar{R}. \end{cases} \end{aligned}$$

- DP algorithm:

$$\bar{J}_k(p_k) = \min \left[(1 - p_k)C, I + \underset{z_{k+1}}{E} \left\{ \bar{J}_{k+1} \left(\Phi(p_k, z_{k+1}) \right) \right\} \right]$$

starting with

$$\bar{J}_{N-1}(p_{N-1}) = \min \left[(1 - p_{N-1})C, I + (1 - t)(1 - p_{N-1})C \right].$$

INSTRUCTION EXAMPLE III

- Write the DP algorithm as

$$\bar{J}_k(p_k) = \min \left[(1 - p_k)C, I + A_k(p_k) \right],$$

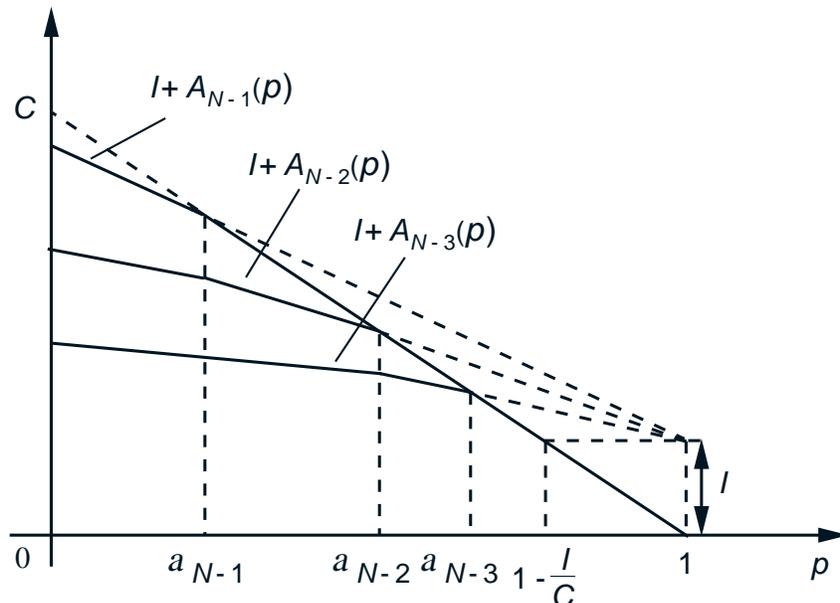
where

$$\begin{aligned} A_k(p_k) = & P(z_{k+1} = R \mid I_k) \bar{J}_{k+1}(\Phi(p_k, R)) \\ & + P(z_{k+1} = \bar{R} \mid I_k) \bar{J}_{k+1}(\Phi(p_k, \bar{R})) \end{aligned}$$

- Can show by induction that $A_k(p)$ are piecewise linear, concave, monotonically decreasing, with

$$A_{k-1}(p) \leq A_k(p) \leq A_{k+1}(p), \quad \text{for all } p \in [0, 1].$$

(The cost-to-go at knowledge prob. p increases as we come closer to the end of horizon.)



MIT OpenCourseWare
<http://ocw.mit.edu>

6.231 Dynamic Programming and Stochastic Control
Fall 2015

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.