

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

6.231 Dynamic Programming and Stochastic Control
Fall 2008

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

6.231 DYNAMIC PROGRAMMING

LECTURE 17

LECTURE OUTLINE

- We start a four-lecture sequence on advanced infinite horizon DP
- We allow infinite state space, so the stochastic shortest path framework cannot be used any more
- The discounted problem is the proper starting point for this analysis
- The central mathematical structure is that the DP mapping is a contraction mapping (instead of existence of a termination state)

DISCOUNTED PROBLEMS W/ BOUNDED COST

- Stationary system with arbitrary state space

$$x_{k+1} = f(x_k, u_k, w_k), \quad k = 0, 1, \dots$$

- Cost of a policy $\pi = \{\mu_0, \mu_1, \dots\}$

$$J_\pi(x_0) = \lim_{N \rightarrow \infty} E_{w_k} \left\{ \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\}$$

with $\alpha < 1$, and for some M , we have $|g(x, u, w)| \leq M$ for all (x, u, w)

- **Shorthand notation for DP mappings** (operate on functions of state to produce other functions)

$$(TJ)(x) = \min_{u \in U(x)} E_w \left\{ g(x, u, w) + \alpha J(f(x, u, w)) \right\}, \quad \forall x$$

TJ is the optimal cost function for the one-stage problem with stage cost g and terminal cost αJ .

- For any stationary policy μ

$$(T_\mu J)(x) = E_w \left\{ g(x, \mu(x), w) + \alpha J(f(x, \mu(x), w)) \right\}, \quad \forall x$$

“SHORTHAND” THEORY – A SUMMARY

- **Cost function expressions** [with $J_0(x) \equiv 0$]

$$J_\pi(x) = \lim_{k \rightarrow \infty} (T_{\mu_0} T_{\mu_1} \cdots T_{\mu_k} J_0)(x), \quad J_\mu(x) = \lim_{k \rightarrow \infty} (T_\mu^k J_0)(x)$$

- **Bellman’s equation:** $J^* = T J^*$, $J_\mu = T_\mu J_\mu$
- **Optimality condition:**

$$\mu: \text{optimal} \quad \langle == \rangle \quad T_\mu J^* = T J^*$$

- **Value iteration:** For any (bounded) J and all x ,

$$J^*(x) = \lim_{k \rightarrow \infty} (T^k J)(x)$$

- **Policy iteration:** Given μ^k ,
 - Policy evaluation: Find J_{μ^k} by solving

$$J_{\mu^k} = T_{\mu^k} J_{\mu^k}$$

- Policy improvement: Find μ^{k+1} such that

$$T_{\mu^{k+1}} J_{\mu^k} = T J_{\mu^k}$$

TWO KEY PROPERTIES

- **Monotonicity property:** For any functions J and J' such that $J(x) \leq J'(x)$ for all x , and any μ

$$(TJ)(x) \leq (TJ')(x), \quad \forall x,$$

$$(T_\mu J)(x) \leq (T_\mu J')(x), \quad \forall x.$$

- **Additivity property:** For any J , any scalar r , and any μ

$$(T(J + re))(x) = (TJ)(x) + \alpha r, \quad \forall x,$$

$$(T_\mu(J + re))(x) = (T_\mu J)(x) + \alpha r, \quad \forall x,$$

where e is the unit function [$e(x) \equiv 1$].

CONVERGENCE OF VALUE ITERATION

- If $J_0 \equiv 0$,

$$J^*(x) = \lim_{N \rightarrow \infty} (T^N J_0)(x), \quad \text{for all } x$$

Proof: For any initial state x_0 , and policy $\pi = \{\mu_0, \mu_1, \dots\}$,

$$\begin{aligned} J_\pi(x_0) &= E \left\{ \sum_{k=0}^{\infty} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \\ &= E \left\{ \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \\ &\quad + E \left\{ \sum_{k=N}^{\infty} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \end{aligned}$$

The tail portion satisfies

$$\left| E \left\{ \sum_{k=N}^{\infty} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \right| \leq \frac{\alpha^N M}{1 - \alpha},$$

where $M \geq |g(x, u, w)|$. Take the min over π of both sides. **Q.E.D.**

BELLMAN'S EQUATION

- The optimal cost function J^* satisfies Bellman's Eq., i.e. $J^* = T(J^*)$.

Proof: For all x and N ,

$$J^*(x) - \frac{\alpha^N M}{1 - \alpha} \leq (T^N J_0)(x) \leq J^*(x) + \frac{\alpha^N M}{1 - \alpha},$$

where $J_0(x) \equiv 0$ and $M \geq |g(x, u, w)|$. Applying T to this relation, and using Monotonicity and Additivity,

$$\begin{aligned} (TJ^*)(x) - \frac{\alpha^{N+1} M}{1 - \alpha} &\leq (T^{N+1} J_0)(x) \\ &\leq (TJ^*)(x) + \frac{\alpha^{N+1} M}{1 - \alpha} \end{aligned}$$

Taking the limit as $N \rightarrow \infty$ and using the fact

$$\lim_{N \rightarrow \infty} (T^{N+1} J_0)(x) = J^*(x)$$

we obtain $J^* = TJ^*$. **Q.E.D.**

THE CONTRACTION PROPERTY

- **Contraction property:** For any bounded functions J and J' , and any μ ,

$$\max_x |(TJ)(x) - (TJ')(x)| \leq \alpha \max_x |J(x) - J'(x)|,$$

$$\max_x |(T_\mu J)(x) - (T_\mu J')(x)| \leq \alpha \max_x |J(x) - J'(x)|.$$

Proof: Denote $c = \max_{x \in S} |J(x) - J'(x)|$. Then

$$J(x) - c \leq J'(x) \leq J(x) + c, \quad \forall x$$

Apply T to both sides, and use the Monotonicity and Additivity properties:

$$(TJ)(x) - \alpha c \leq (TJ')(x) \leq (TJ)(x) + \alpha c, \quad \forall x$$

Hence

$$|(TJ)(x) - (TJ')(x)| \leq \alpha c, \quad \forall x.$$

Q.E.D.

IMPLICATIONS OF CONTRACTION PROPERTY

- Bellman's equation $J = TJ$ has a unique solution, namely J^* , and for any bounded J , we have

$$\lim_{k \rightarrow \infty} (T^k J)(x) = J^*(x), \quad \forall x$$

Proof: Use

$$\begin{aligned} \max_x |(T^k J)(x) - J^*(x)| &\leq \max_x |(T^k J)(x) - (T^k J^*)(x)| \\ &\leq \alpha^k \max_x |J(x) - J^*(x)| \end{aligned}$$

- **Convergence rate:** For all k ,

$$\max_x |(T^k J)(x) - J^*(x)| \leq \alpha^k \max_x |J(x) - J^*(x)|$$

- Also, for each stationary μ , J_μ is the unique solution of $J = T_\mu J$ and

$$\lim_{k \rightarrow \infty} (T_\mu^k J)(x) = J_\mu(x), \quad \forall x,$$

for any bounded J .

NEC. AND SUFFICIENT OPT. CONDITION

- A stationary policy μ is optimal if and only if $\mu(x)$ attains the minimum in Bellman's equation for each x ; i.e.,

$$TJ^* = T_\mu J^*.$$

Proof: If $TJ^* = T_\mu J^*$, then using Bellman's equation ($J^* = TJ^*$), we have

$$J^* = T_\mu J^*,$$

so by uniqueness of the fixed point of T_μ , we obtain $J^* = J_\mu$; i.e., μ is optimal.

- Conversely, if the stationary policy μ is optimal, we have $J^* = J_\mu$, so

$$J^* = T_\mu J^*.$$

Combining this with Bellman's equation ($J^* = TJ^*$), we obtain $TJ^* = T_\mu J^*$. **Q.E.D.**

COMPUTATIONAL METHODS

- **Value iteration** and variants
 - Gauss-Seidel version
 - Approximate value iteration
- **Policy iteration** and variants
 - Combination with value iteration
 - Modified policy iteration
 - Asynchronous policy iteration
- **Linear programming**

$$\text{maximize } \sum_{i=1}^n J(i)$$

$$\text{subject to } J(i) \leq g(i, u) + \alpha \sum_{j=1}^n p_{ij}(u) J(j), \quad \forall (i, u)$$

- **Approximate linear programming:** use in place of $J(i)$ a low-dim. basis function representation

$$\tilde{J}(i, r) = \sum_{k=1}^m r_k w_k(i)$$

and low-dim. LP (with many constraints)