

# CE 273

## Markov Decision Processes

Lecture 16

### **Policy Iteration and Linear Programming for Average Cost MDPs**

# Previously on Markov Decision Processes

## Proposition

*For any transition matrix  $P$  and fundamental matrix  $H$*

$$P^* = PP^* = P^*P = P^*P^*$$

$$P^*H = HP^* = 0$$

$$P^* + H = I + PH$$

# Previously on Markov Decision Processes

## Definition

The gain  $J_\mu$  of a policy  $\mu$  is defined as

$$J_\mu = P_\mu^* g_\mu$$

## Definition

The bias  $h_\mu$  of a policy  $\mu$  is defined as

$$h_\mu = H_\mu g_\mu$$

where  $H_\mu = (I - P_\mu + P_\mu^*)^{-1} - P_\mu^*$  and is called the fundamental matrix.

In addition, suppose the associated Markov chain is aperiodic, i.e., if  $P_\mu^* = \lim_{N \rightarrow \infty} P_\mu^N$  (Case III), then we can interpret  $h_\mu$  as

$$h_\mu = \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} P_\mu^k (g_\mu - J_\mu)$$

a relative cost vector, i.e., the difference of the total cost of  $\mu$  and the total cost if one-stage costs were set to  $J_\mu$ .

# Previously on Markov Decision Processes

Unlike discounted and total cost MDPs, where we could solve a system of equations for a given policy (and use this in the policy iteration algorithm), we cannot simply solve

$$\begin{aligned}J &= P_\mu J \\ J + h &= g_\mu + P_\mu h\end{aligned}$$

to get the average cost of policy  $\mu$ . (Why?)

If  $(J_\mu, h_\mu)$  solves the above system, then  $(J_\mu, h_\mu + \text{constant})$  also satisfies the above system. Hence, there are an infinite number of solutions. We will call these policy evaluation equations for easy referencing.

In general, it can be shown that all solutions to the above system are of the form  $(J_\mu, h_\mu + d)$ , where  $d = P_\mu d$ .

# Previously on Markov Decision Processes

The earlier proposition and discussion established that a Blackwell optimal policy is optimal to the average cost problem.

Further, optimal policies were found to satisfy some equations which are the necessary conditions for optimality. It can also be shown that they are sufficient.

## Proposition

*If  $J'$  and  $h'$  satisfy the following pair of optimality equations*

$$J(i) = \min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) J(j) \quad \forall i = 1, \dots, n$$

$$J(i) + h(i) = \min_{u \in \bar{U}(i)} \left\{ g(i, u) + \sum_{j=1}^n p_{ij}(u) h(j) \right\} \quad \forall i = 1, \dots, n$$

*where  $\bar{U}(i)$  is the set of controls that attain the minimum in the above equation. Then,  $J' = J^*$  is the optimal average cost vector.*

*Further, if a stationary policy  $\mu$  attains the minimum in the above equations, then it is the optimal policy  $\mu^*$ .*

# Previously on Markov Decision Processes

In summary, if the average cost is independent of the initial state, the following proposition is true

## Proposition

*If a scalar  $\lambda$  and a vector  $h$  satisfy*

$$\lambda + h(i) = \min_{u \in U(i)} \left\{ g(i, u) + \sum_{j=1}^n p_{ij}(u)h(j) \right\} \quad \forall i = 1, \dots, n$$

*then  $\lambda$  is the optimal average cost  $J^*(i)$  for all  $i$ , i.e.,*

$$\lambda = \min_{\mu} J_{\mu}(i) = J^*(i) \quad \forall i = 1, \dots, n$$

*Further, if  $\mu^*$  attains the minimum in the first expression, then  $J_{\mu^*}(i) = \lambda \forall i$ .*

In shorthand, the first equation can be rewritten as  $\lambda e + h = Th$ . Think of this as being analogous to  $J^* = TJ^*$  in the discounted world.

# Previously on Markov Decision Processes

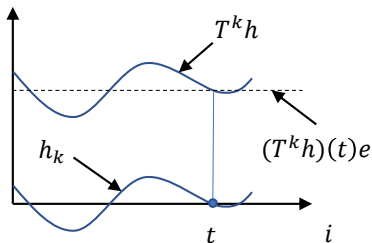
## Theorem

$$J^* = \lim_{k \rightarrow \infty} \frac{1}{k} T^k h$$

For unichain MDPs,

$$h_k = T^k h - (T^k h)(t)e$$

where  $t$  is some arbitrary state. Effectively, we are shifting the entire function by a constant. But note that the constant varies across iterations.



The iterates  $h_k$  remain bounded and the bounds do not depend on  $k$ .

# Previously on Markov Decision Processes

The relative value iteration (RVI) works for unichain MDPs which induce aperiodic Markov Chains.

Define the span semi-norm of a vector  $h$  as

$$sp(h) = \max_{i \in X} h(i) - \min_{i \in X} h(i)$$

---

## RELATIVE VALUE ITERATION

---

Fix a tolerance level  $\epsilon > 0$  and select a state  $t$

Select  $h_0 \in B(X)$  and  $k \leftarrow 0$

$h_1 \leftarrow Th_0 - (Th_0)(t)e$

**while**  $sp(h_{k+1} - h_k) > \epsilon$  **do**

$k \leftarrow k + 1$

$h_{k+1} \leftarrow Th_k - (Th_k)(t)e$

**end while**

Select  $\mu$  such that

$$\mu(i) \in \arg \min_{u \in U(i)} \left\{ g(i, u) + \sum_{j=1}^n p_{ij}(u) h_k(j) \right\}$$



# Lecture Outline

- 1 Policy Iteration
- 2 Linear Programming

## Policy Iteration

# Policy Iteration

## Introduction

A policy iteration algorithm that alternates between policy evaluation and policy improvement can be used to solve the average cost problem.

Policy iteration works for both unichain and multichain MDPs but the steps for the latter type of problem are more involved.

# Policy Iteration

## Policy Evaluation for Unichain MDPs

Consider a unichain MDP. The MDP associated with any policy thus has a single closed communicating class and a set of transient states.

Suppose at the  $k$ th iteration, we have a policy  $\mu_k$ . Then, the policy evaluation is done by solving the system

$$\lambda_k e + h_k = T_{\mu_k} h_k$$

Note that for unichain MDPs  $J_k = P_{\mu_k} J_k$  is always satisfied. However, as noted earlier, the above system does not have a unique solution.

Hence, we select an arbitrary state  $t$  and set  $h_k(t) = 0$ . It can be shown that this new system

$$\begin{aligned}\lambda_k e + h_k &= T_{\mu_k} h_k \\ h_k(t) &= 0\end{aligned}$$

has a unique solution and  $\lambda_k$  corresponds to the gain of the policy  $\mu_k$ .

# Policy Iteration

## Policy Improvement for Unichain MDPs

Note that a  $h_k$  that satisfies the above set of equations need not equal  $h_{\mu_k}$ , which is why we use the subscript  $k$  and not  $\mu_k$ . (We can however call  $\lambda_k$  as  $\lambda_{\mu_k}$ .)

Since we know the policy, we could as well compute the gain  $J_{\mu_k} = \lambda_k e$  and  $h_{\mu_k}$  and use it in the next step but it would require more computation.

Policy improvement is done by finding the controls which optimize  $Th_k$ , that is,

$$T_{\mu_{k+1}} h_k = Th_k$$

The new policy  $\mu_{k+1}$  is the same for any  $h_k$  that satisfies the policy evaluation equation. (Why?)

As before,

- ▶ The algorithm is terminated when  $\mu_{k+1} = \mu_k$ .
- ▶ Ties are broken such that  $\mu_{k+1}(i) = \mu_k(i)$  whenever possible.

# Policy Iteration

## Pseudocode

---

### POLICY ITERATION

---

$k \leftarrow 0$

Pick an initial policy  $\mu_0$  (say a Greedy policy) and some state  $t$

**do**

    Compute  $\lambda_k$  and  $h_k$  by solving i.e.,

▷ **Policy Evaluation**

$$\lambda_k \mathbf{e} + h_k = T_{\mu_k} h_k$$

$$h_k(t) = 0$$

    Compute a new policy  $\mu_{k+1}$  that satisfies

▷ **Policy Improvement**

$$T_{\mu_{k+1}} h_k = T h_k$$

$k \leftarrow k + 1$

**while**  $\mu_{k+1} \neq \mu_k$

$\mu^* \leftarrow \mu_k$  and  $J^* \leftarrow \lambda_k \mathbf{e}$

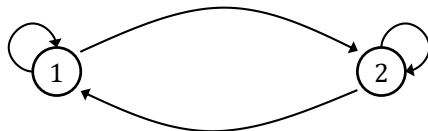
---

Set  $\mu_{k+1}(i) = \mu_k(i)$  whenever possible.

# Policy Iteration

## Example

Perform two iterations of the PI algorithm for the following example with two states 1 and 2. Assume that state 1 is the reference state  $t$ . Start with the policy  $\mu_0(1) = u_1$  and  $\mu_0(2) = u_2$ .



- ▶  $U(1) = \{u_1, u_2\}$
- ▶  $g(1, u_1) = 2, g(1, u_2) = 0.5$
- ▶  $p_{1j}(u_1) = [3/4 \ 1/4]$
- ▶  $p_{1j}(u_2) = [1/4 \ 3/4]$
- ▶  $U(2) = \{u_1, u_2\}$
- ▶  $g(2, u_1) = 1, g(2, u_2) = 3$
- ▶  $p_{2j}(u_1) = [3/4 \ 1/4]$
- ▶  $p_{2j}(u_2) = [1/4 \ 3/4]$

# Policy Iteration

## Main Result

Let's now see why the policy iteration method works.

### Proposition

Consider an unichain MDP and a policy  $\mu$  with gain-bias pair  $(\lambda_\mu, h_\mu)$ . Suppose  $\mu'$  is obtained from  $T_{\mu'} h_\mu = Th_\mu$  and denote using  $(\lambda_{\mu'}, h_{\mu'})$  the gain-bias pair of  $\mu'$ . If  $\mu' \neq \mu$ , the one of the following is true

- 1  $\lambda_{\mu'} < \lambda_\mu$
- 2  $\lambda_{\mu'} = \lambda_\mu$  and  $h_{\mu'}(i) \leq h_\mu(i)$  for all  $i = 1, \dots, n$  with equality occurring for states that are recurrent and strict inequality for at least one transient state.

### Proof.

We will only prove  $\lambda_{\mu'} \leq \lambda_\mu$ . To do so, it is enough to show

$$P_{\mu'}^*(T_\mu h_\mu - T_{\mu'} h_\mu) = (\lambda_\mu - \lambda_{\mu'})e$$

(Why?) By construction of  $\mu'$ ,  $T_{\mu'} h_\mu = Th_\mu \leq T_\mu h_\mu$ . Therefore,  $P_{\mu'}^*(T_\mu h_\mu - T_{\mu'} h_\mu) \geq 0$ .



# Policy Iteration

## Main Result

Proof.

$$P_{\mu'}^*(T_{\mu}h_{\mu} - T_{\mu'}h_{\mu}) = (\lambda_{\mu} - \lambda_{\mu'})e$$

Consider the LHS:

$$\begin{aligned} P_{\mu'}^*(T_{\mu}h_{\mu} - T_{\mu'}h_{\mu}) &= P_{\mu'}^*\left(T_{\mu}h_{\mu} - \left(T_{\mu'}h_{\mu'} + (T_{\mu'}h_{\mu} - T_{\mu'}h_{\mu'})\right)\right) \\ &= P_{\mu'}^*\left(\lambda_{\mu}e + h_{\mu} - \left(\lambda_{\mu'}e + h_{\mu'} + (T_{\mu'}h_{\mu} - T_{\mu'}h_{\mu'})\right)\right) \\ &= P_{\mu'}^*\left(\lambda_{\mu}e + h_{\mu} - \left(\lambda_{\mu'}e + h_{\mu'} + P_{\mu'}(h_{\mu} - h_{\mu'})\right)\right) \\ &= P_{\mu'}^*\left((\lambda_{\mu} - \lambda_{\mu'})e + (I - P_{\mu'})(h_{\mu} - h_{\mu'})\right) \\ &= P_{\mu'}^*(\lambda_{\mu} - \lambda_{\mu'})e + (P_{\mu'}^* - P_{\mu'}^*P_{\mu'})(h_{\mu} - h_{\mu'}) \\ &= (\lambda_{\mu} - \lambda_{\mu'})e + 0 \end{aligned}$$



# Policy Iteration

## Policy Evaluation for Multichain MDPs

For multichain MDPs, the same ideas work but both policy evaluation and improvement involve more equations.

Recall that the gain and bias of a policy satisfy

$$\begin{aligned}J &= P_\mu J \\ J + h &= g_\mu + P_\mu h\end{aligned}$$

One cannot solve this system and identify the bias since it is not unique.

Non-uniqueness was an issue even in unichain MDPs, but the choice of the bias did not matter when we perform policy improvement.

However, that is no longer true for multichain MDPs. At every iteration  $k$ , we need  $(J_{\mu_k}, h_{\mu_k})$  to find an improved policy  $\mu_{k+1}$ !

# Policy Iteration

## Policy Evaluation for Multichain MDPs

An obvious option to find  $(J_{\mu_k}, h_{\mu_k})$  is to estimate  $P_{\mu_k}^*$  and the fundamental matrix  $H_{\mu_k}$ . But this is computationally expensive. Alternately, the following result can be used

### Proposition

*Consider a stationary policy  $\mu$  with the gain-bias pair  $(J_\mu, h_\mu)$ . The set of solutions  $(J, h, v)$  to the following equations*

$$J = P_\mu J$$

$$J + h = g_\mu + P_\mu h$$

$$h + v = P_\mu v$$

*are of the form  $(J_\mu, h_\mu, -H_\mu^2 g_\mu + d)$  where  $d$  satisfies  $d = P_\mu d$ .*

# Policy Iteration

## Policy Improvement for Multichain MDPs

Once we have  $(J_{\mu_k}, h_{\mu_k})$ ,  $\mu_{k+1}$  is obtained from the following two-stage policy improvement procedure:

### Step 1:

Choose a policy  $\mu_{k+1}$  which satisfies

$$P_{\mu_{k+1}} J_{\mu_k} = \min_{\mu} P_{\mu} J_{\mu_k}$$

In other words,

$$\mu_{k+1}(i) \in \arg \min_{u \in U(i)} \left\{ \sum_{j=1}^n p_{ij}(u) J_{\mu_k}(j) \right\}$$

while setting  $\mu_{k+1}(i) = \mu_k(i)$  whenever possible. If  $\mu_{k+1} \neq \mu_k$ , then we can switch to the policy evaluation procedure. Else, go to Step 2.

# Policy Iteration

## Policy Improvement for Multichain MDPs

### Step 2:

Choose a policy  $\mu_{k+1}$  which satisfies

$$P_{\mu_{k+1}} J_{\mu_k} = \min_{\mu} P_{\mu} J_{\mu_k}$$

$$T_{\mu_{k+1}} h_{\mu_k} = \min_{\mu \in \bar{\Pi}} T_{\mu} h_{\mu_k}$$

where  $\bar{\Pi}$  is the set of policies which attain the minimum in  $\min_{\mu} P_{\mu} J_{\mu_k}$ . Alternately we can write,

$$\mu_{k+1}(i) \in \arg \min_{u \in U(i)} \left\{ \sum_{j=1}^n p_{ij}(u) J_{\mu_k}(j) \right\}$$

$$\mu_{k+1}(i) \in \arg \min_{u \in \bar{U}(i)} \left\{ g(i, u) + \sum_{j=1}^n p_{ij}(u) h_{\mu_k}(j) \right\}$$

where  $\bar{U}(i)$  is the set of controls are the optima of  $\sum_{j=1}^n p_{ij}(u) J_{\mu_k}(j)$ . Again set  $\mu_{k+1}(i) = \mu_k(i)$  whenever possible.

# Policy Iteration

## Main Result

The above policy iteration procedure for multichain MDPs works because of the following proposition (which is very similar to what we saw in unichain MDPs).

### Proposition

Let  $\mu_k$  be a policy with gain-bias pair  $(J_{\mu_k}, h_{\mu_k})$ . Suppose that  $\mu_{k+1}$  is obtained from policy improvement step and let  $(J_{\mu_{k+1}}, h_{\mu_{k+1}})$  be its gain-bias pair. If  $\mu_{k+1} \neq \mu_k$  then one of the following is true

- 1  $J_{\mu_{k+1}}(i) \leq J_{\mu_k}(i)$  for all  $i = 1, \dots, n$  with a strict inequality for at least one state  $i$ .
- 2  $J_{\mu_{k+1}} = J_{\mu_k}$  and  $h_{\mu_{k+1}} \leq h_{\mu_k}$  for all  $i = 1, \dots, n$  with strict inequality for at least one state  $i$  transient under  $\mu_{k+1}$ .

## Linear Programming

# Linear Programming

## Introduction

Both unichain and multichain MDPs can also be solved using linear programming.

We will need the following proposition to set up the LPs

### Proposition

*Let  $J$  and  $h$  be vectors which satisfy*

$$\begin{aligned} J &\leq P_\mu J \\ J + h &\leq T_\mu h \end{aligned}$$

*Then,  $J \leq J_\mu$ . Further, if equality holds in the first two inequalities,  $J = J_\mu$ .*



# Linear Programming

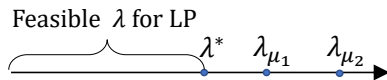
## Unichain MDPs

Consider the case of unichain MDPs. The optimality conditions can be written as

$$\lambda e + h = Th = \min_{\mu} T_{\mu}h$$

Alternately, we can write  $\lambda e + h \leq T_{\mu}h$  for all stationary policies  $\mu$ . And by setting  $J = \lambda e$ , the previous proposition can be used.

Thus, for all functions  $h$  and scalars  $\lambda$  that satisfy  $\lambda e + h \leq T_{\mu}h$ ,  $\lambda \leq \lambda_{\mu}$  for every stationary policy  $\mu$ .



Further for the optimal policy  $\mu^*$ , the optimal average cost  $\lambda^*$  satisfies  $\lambda^* e + h = T_{\mu^*}h$ . Thus,  $\lambda^*$  is the largest  $\lambda$  that satisfies  $\lambda e + h \leq Th$ .

# Linear Programming

## Primal Problem

The primal LP for unichain MDPs can thus be written as

$$\begin{aligned} & \max \lambda \\ \text{s.t. } & \lambda + h(i) \leq g(i, u) + \sum_{j=1}^n p_{ij}(u)h(j) \quad \forall i = 1, \dots, n, u \in U(i) \end{aligned}$$

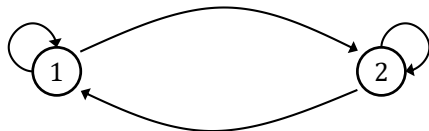
which in standard form looks like

$$\begin{aligned} & \max \lambda \\ \text{s.t. } & \lambda + h(i) - \sum_{j=1}^n p_{ij}(u)h(j) \leq g(i, u) \quad \forall i = 1, \dots, n, u \in U(i) \end{aligned}$$

# Linear Programming

## Example

Write the primal LP for the following example with two states 1 and 2.



- ▶  $U(1) = \{u_1, u_2\}$
- ▶  $g(1, u_1) = 2, g(1, u_2) = 0.5$
- ▶  $p_{1j}(u_1) = [3/4 \ 1/4]$
- ▶  $p_{1j}(u_2) = [1/4 \ 3/4]$
- ▶  $U(2) = \{u_1, u_2\}$
- ▶  $g(2, u_1) = 1, g(2, u_2) = 3$
- ▶  $p_{2j}(u_1) = [3/4 \ 1/4]$
- ▶  $p_{2j}(u_2) = [1/4 \ 3/4]$

# Linear Programming

## Dual Problem

Write the dual of the above LP.

- ▶ The number of dual variables equal to the number of constraints in the primal.
- ▶ Since the primal constraints are of the  $\leq$  form, the dual variables must be  $\geq 0$ .
- ▶ Since the primal variables are unconstrained, the dual will have equality constraints. (How many?)

# Linear Programming

## Dual Problem

Define variables  $z(i, u)$  where  $i \in X, u \in U(i)$ .

$$\min \sum_{i=1}^n \sum_{u \in U(i)} g(i, u) z(i, u)$$

$$\text{s.t.} \quad \sum_{u \in U(i)} z(i, u) - \sum_{j=1}^n \sum_{u \in U(j)} p_{ji}(u) z(j, u) = 0 \quad \forall i = 1, \dots, n$$

$$\sum_{i=1}^n \sum_{u \in U(i)} z(i, u) = 1$$

$$z(i, u) \geq 0 \quad \forall u \in U(i), i = 1, \dots, n$$

If you think of  $\sum_{u \in U(i)} z(i, u)$  as a new variable  $z_i$ , what do the constraints represent?

# Linear Programming

## Constructing Solutions from LPs

The primal problem gives us only the optimal  $\lambda$ . The dual solution on the other hand can be used to construct the optimal policy as well. (How did we do this for the discounted problem?)

Let  $z^*(i, u)$  be the optimal dual solution. Define for all  $i$ ,

$$U^*(i) = \{u \in U(i) \mid z^*(i, u) > 0\}$$

For average cost MDPs, it is not necessary that the above set is a singleton. Define a new set

$$C^* = \left\{ i \mid \sum_{u \in U(i)} z^*(i, u) > 0 \right\}$$

Note that the sets  $U^*(i)$  and  $C^*$  are non-empty. (Why?) The following policy can be shown to be optimal

$$\mu^*(i) = \begin{cases} \text{any } u \in U^*(i) & \text{if } i \in C^* \\ \text{any } u \in U(i) & \text{if } i \notin C^* \end{cases}$$

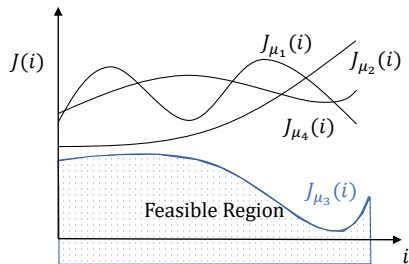
It can also be shown that that the set  $C^*$  is a closed communicating class under the optimal policy  $\mu^*$ .

# Linear Programming

## Multichain MDPs

The earlier proposition can also be used to construct a LP for the multichain MDP. Recall that we now have two sets of feasible constraints  $J \leq P_\mu J$  and  $J + h \leq T_\mu h$ .

This implies that any pair of vectors  $(J, h)$  that satisfies the above system,  $J \leq J_\mu$ .



Thus,  $J^*$  is the “largest” vector satisfying the two constraints.

# Linear Programming

## Primal LP for Multichain MDPs

The primal problem for the multichain MDP can be written as

$$\begin{aligned} & \max \sum_{i=1}^n a_i y(i) \\ \text{s.t. } & y(i) \leq \sum_{j=1}^n p_{ij}(u) y(j) \quad \forall i = 1, \dots, n, u \in U(i) \\ & y(i) + h(i) \leq g(i, u) + \sum_{j=1}^n p_{ij}(u) h(j) \quad \forall i = 1, \dots, n, u \in U(i) \end{aligned}$$

where  $a$  is a row vector of strictly positive reals satisfying  $\sum_i a_i = 1$ .



# Linear Programming

## Dual LP for Multichain MDPs

Hence, the dual takes the form

$$\min \sum_{i=1}^n \sum_{u \in U(i)} g(i, u) z(i, u)$$

$$\text{s.t.} \quad \sum_{u \in U(i)} z(i, u) - \sum_{j=1}^n \sum_{u \in U(j)} p_{ji}(u) z(j, u) = 0 \quad \forall i = 1, \dots, n$$

$$\sum_{u \in U(i)} (z(i, u) + r(i, u)) - \sum_{j=1}^n \sum_{u \in U(j)} r(j, u) p_{ji}(u) = a_i \quad \forall i = 1, \dots, n$$

$$z(i, u) \geq 0 \quad \forall u \in U(i), i = 1, \dots, n$$

$$r(i, u) \geq 0 \quad \forall u \in U(i), i = 1, \dots, n$$

# Linear Programming

## Dual LP for Multichain MDPs

As before, the optimal policy is constructed by first supposing

$$C^* = \left\{ i \mid \sum_{u \in U(i)} z^*(i, u) > 0 \right\}$$

using which, we define

$$\mu^*(i) = \begin{cases} \text{any } u \text{ such that } z^*(i, u) > 0 & \text{if } i \in C^* \\ \text{any } u \text{ such that } r^*(i, u) > 0 & \text{if } i \notin C^* \end{cases}$$

# Your Moment of Zen

