

Dynamic Programming Lecture #8

Outline:

- Controlled Markov chains
- Value iteration proof

Controlled Markov Chains

- State transition probabilities are now a function of control input: $P(u)$
- Example: Machine repair.

- States = $\begin{cases} 1 & \text{up} \\ 2 & \text{down} \end{cases}$
- Controls = $\begin{cases} S & \text{stop \& repair if needed} \\ C & \text{continue} \end{cases}$
- Control dependent transitions:

$$P(C) = \begin{pmatrix} Pr(\text{up}|\text{up}) & Pr(\text{down}|\text{up}) \\ Pr(\text{up}|\text{down}) & Pr(\text{down}|\text{down}) \end{pmatrix} = \begin{pmatrix} 0.8 & 0.2 \\ 0 & 1 \end{pmatrix}$$

- For a policy, $u = \mu(x)$,

$$P(\mu) = \begin{pmatrix} \text{—1st row of } P(\mu(1))\text{—} \\ \text{—2nd row of } P(\mu(2))\text{—} \\ \vdots \\ \text{—kth row of } P(\mu(k))\text{—} \\ \vdots \\ \text{—n-th row of } P(\mu(n))\text{—} \end{pmatrix}$$

Controlled Markov Chains & Random Disturbances

- Original (stage invariant) setup:

$$x_{k+1} = f(x_k, u_k, w_k)$$

- Controlled Markov chains:

$$\text{P}[x_{k+1} = j | x_k = i, u_k] = p_{ij}(u_k)$$

What happened to w ?

- Connection:

$$x_{k+1} = w_k$$

$$\text{P}[w_k = j | x_k = i, u_k] = p_{ij}(u_k)$$

- Value iteration:

- Assume stage cost does not depend on w (no loss of generality)
- Original setup:

$$J_k(x_k) = \min_{u_k \in \mathcal{U}_k(x_k)} \text{E}[g_k(x_k, u_k) + J_{k+1}(f(x_k, u_k, w_k))]$$

- Markov chain setup:

$$J_k(i) = \min_{u_k \in \mathcal{U}_k(i)} g_k(i, u_k) + \sum_{j=1}^n J_{k+1}(j) p_{ij}(u_k)$$

- Can also have *stage dependent* probabilities

Cost of Policy

- Policy: $\pi = \{\mu_1, \mu_2, \dots, \mu_{N-1}\}$
- Transitions: $P_k(\mu_k)$
- i^{th} coordinate row vector:

$$e_i = (0 \quad \dots \quad 0 \quad 1 \quad 0 \quad \dots \quad 0)$$

- Policy dependent cost vector:

$$G_k(\mu) = \begin{pmatrix} g_k(1, \mu(1)) \\ g_k(2, \mu(2)) \\ \vdots \\ g_k(n, \mu(n)) \end{pmatrix} \quad \& \quad G_N = \begin{pmatrix} g_N(1) \\ g_N(2) \\ \vdots \\ g_N(n) \end{pmatrix}$$

- Cost of trajectory $\{x_0, x_1, \dots, x_N\}$ under π :

$$J(x_0, \dots, x_N; \pi) = e_{x_0} G_0(\mu_0) + e_{x_1} G_1(\mu_1) + \dots + e_{x_N} G_N$$

- *Expected* cost of trajectory given x_0 under π :

$$\sum_{x_0 \dots x_N} J(x_0, \dots, x_{N-1}; \pi) P[x_0, \dots, x_N | x_0, \pi]$$

Cost of Policy, cont

- Use $\mathbb{E} [\sum_j Z_j] = \sum_j \mathbb{E} [Z_j]$ on

$$J(x_0, \dots, x_N; \pi) = e_{x_0} G_0(\mu_0) + e_{x_1} G_1(\mu_1) + \dots + e_{x_N} G_N$$

to rewrite *expected cost* given x_0 under π :

$$\begin{aligned} \mathbb{E} [J(x_0, \dots, x_{N-1}; \pi) | x_0] &= e_{x_0} G_0(\mu_0) \\ &\quad + \underbrace{e_{x_0} P_0(\mu_0)}_{q_1 = \mathbb{E}[e_{x_1} | x_0, \mu_0]} G_1(\mu_1) \\ &\quad + \underbrace{e_{x_0} P_0(\mu_0) P_1(\mu_1)}_{q_2 = \mathbb{E}[e_{x_2} | x_0, \mu_0, \mu_1]} G_2(\mu_2) + \dots \\ &\quad + \underbrace{e_{x_0} P_0(\mu_0) P_1(\mu_1) \dots P_{N-2}(\mu_{N-2})}_{q_{N-1} = \mathbb{E}[e_{x_{N-1}} | x_0, \mu_0, \mu_1, \dots, \mu_{N-2}]} G_{N-1}(\mu_{N-1}) \\ &\quad + \underbrace{e_{x_0} P_0(\mu_0) P_1(\mu_1) \dots P_{N-2}(\mu_{N-2}) P_{N-1}(\mu_{N-1})}_{q_N = \mathbb{E}[e_{x_N} | x_0, \mu_0, \mu_1, \dots, \mu_{N-1}]} G_N \end{aligned}$$

- Notice triangular structure in μ_k dependence!

Value Iteration Revisited

- Rewrite (for some function α and row vector β):

$$\begin{aligned} \mathbb{E} [J(x_0, \dots, x_{N-1}; \pi) | x_0] = \\ \alpha(x_0, \mu_0, \dots, \mu_{N-2}) + \beta(x_0, \mu_0, \dots, \mu_{N-2}) (G_{N-1}(\mu_{N-1}) + P_{N-1}(\mu_{N-1})G_N) \end{aligned}$$

- Accordingly:

$$\begin{aligned} \min_{\mu_0, \dots, \mu_{N-1}} \mathbb{E} [J(x_0, \dots, x_{N-1}; \pi) | x_0] = \\ \min_{\mu_0, \dots, \mu_{N-1}} \alpha(x_0, \mu_0, \dots, \mu_{N-2}) \\ + \beta(x_0, \mu_0, \dots, \mu_{N-2}) \begin{pmatrix} g_{N-1}(1, \mu_{N-1}(1)) + \sum_{j=1}^n P_{N-1}^{1j}(\mu_{N-1}(1))G_N(j) \\ g_{N-1}(2, \mu_{N-1}(2)) + \sum_{j=1}^n P_{N-1}^{2j}(\mu_{N-1}(2))G_N(j) \\ \vdots \\ g_{N-1}(n, \mu_{N-1}(n)) + \sum_{j=1}^n P_{N-1}^{nj}(\mu_{N-1}(n))G_N(j) \end{pmatrix} \end{aligned}$$

- Minimizing over μ_{N-1} results in term-by-term minimization:

$$J_{N-1}^* = \begin{pmatrix} \min_{u_{N-1}} g_{N-1}(1, u_{N-1}) + \sum_{j=1}^n P_{N-1}^{1j}(u_{N-1})G_N(j) \\ \min_{u_{N-1}} g_{N-1}(2, u_{N-1}) + \sum_{j=1}^n P_{N-1}^{2j}(u_{N-1})G_N(j) \\ \vdots \\ \min_{u_{N-1}} g_{N-1}(n, u_{N-1}) + \sum_{j=1}^n P_{N-1}^{nj}(u_{N-1})G_N(j) \end{pmatrix}$$

Value Iteration Revisited, cont

- μ_{N-1} eliminated (replaced by J_{N-1}^*):

$$\begin{aligned} \min_{\mu_0, \dots, \mu_{N-1}} \mathbb{E} [J(x_0, \dots, x_{N-1}; \pi) | x_0] &= \\ \min_{\mu_0, \dots, \mu_{N-2}} \alpha(x_0, \mu_0, \dots, \mu_{N-2}) + \beta(x_0, \mu_0, \dots, \mu_{N-2}) J_{N-1}^* \end{aligned}$$

- Repeat analysis, but now J_{N-1}^* plays role of G_N :

$$\begin{aligned} \min_{\mu_0, \dots, \mu_{N-1}} \mathbb{E} [J(x_0, \dots, x_{N-1}; \pi) | x_0] &= \\ \min_{\mu_0, \dots, \mu_{N-2}} \alpha(x_0, \mu_0, \dots, \mu_{N-3}) + \beta(x_0, \mu_0, \dots, \mu_{N-3}) (G_{N-2}(\mu_{N-2}) + P_{N-2}(\mu_{N-2}) J_{N-1}^*) \end{aligned}$$

- Minimization over μ_{N-2} produces J_{N-2}^* as function of J_{N-1}^*

- Repeated application is same as value iteration:

$$\begin{aligned} J_k(x_k) &= \min_{u_k} \mathbb{E} [g_k(x_k, u_k) + J_{k+1}(f_k(x_k, u_k, w_k))] \\ &= \min_{u_k} g_k(x_k, u_k) + \sum_{j=1}^n J_{k+1}(j) \mathbb{P}[x_{k+1} = j | x_k, u_k] \end{aligned}$$