

Planning in MDPs

Wen Sun

CS 6789: Foundations of Reinforcement Learning

Announcements

HW0 is due Feb 1st

Recap: Value iteration

$$Q^{t+1} = \mathcal{T} Q^t$$

Recap: Value iteration

$$Q^{t+1} = \mathcal{T} Q^t$$

$$\pi^t : \pi^t(s) = \arg \max_a Q^t(s, a)$$

Theorem: $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

Recap: Value iteration

$$Q^{t+1} = \mathcal{T} Q^t$$

$$\pi^t : \pi^t(s) = \arg \max_a Q^t(s, a)$$

Theorem: $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

Q: when will π^t be the optimal policy?

Outline

1. Policy Iteration

2. Computation complexity of VI and PI

3. Linear Programming formulation

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^t}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s, a)} Q^{\pi^t}(s', \pi(s'))$$

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$


3. Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Policy Iteration Algorithm:

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

3. Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$



Policy Iteration Algorithm:

Closed-form for PE
(see 1.1.3 in Monograph)

1. Initialization: $\pi^0 : S \mapsto A$

2. Policy Evaluation: $Q^{\pi^t}(s, a), \forall s, a$

3. Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$



Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

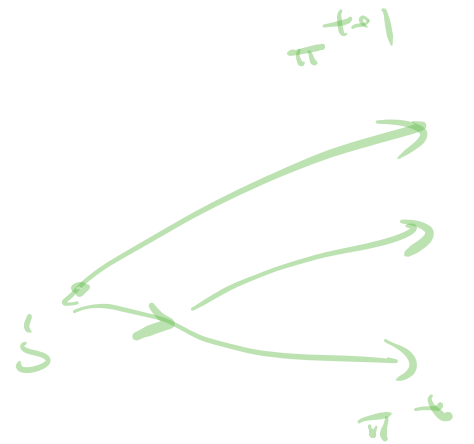
Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - \underline{Q^{\pi^t}(s, a)} = \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

$$= \gamma \mathbb{E}_{s' \sim P(s, a)} \left(Q^{\pi^t}(s', \pi^t(s')) \right) \leftarrow \text{Bell-Equ}$$



Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned} Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[\underbrace{Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s'))}_{\geq 0} + \underbrace{Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s'))}_{\geq 0} \right] \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned}
 Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \\
 &\geq \gamma \mathbb{E}_{s \sim P(s, a)} \left[\gamma \mathbb{E}_{s'' \sim P(s', \pi^{t+1}(s'))} \left[Q^{\pi^{t+1}}(s'', \pi^{t+1}(s'')) - Q^{\pi^t}(s'', \pi^t(s'')) \right] \right]
 \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned} Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[\underbrace{Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s'))}_{\in \left[-\frac{1}{1-\gamma}, \frac{1}{1-\gamma} \right]} \right] \geq \dots, \geq -\gamma^\infty / (1 - \gamma) = 0 \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned} Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \geq \dots, \geq -\gamma^\infty / (1 - \gamma) = 0 \end{aligned}$$

$$\begin{aligned} V^{\pi^{t+1}}(s) &\geq V^{\pi^t}(s), \forall s \\ \left(\overset{t+1}{Q}(s, \pi^{t+1}(s)) \right) &\geq \left(\overset{\pi^t}{Q}(s, \pi^{t+1}(s)) \right) \\ &\geq \left(\overset{\pi^t}{Q}(s, \pi^t(s)) \right) \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

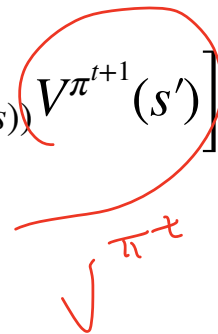
Theorem: Convergence $\|V^{\pi^{t+1}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^t} - V^{\star}\|_{\infty}$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^t} - V^{\star}\|_{\infty}$

$$V^{\star}(s) - V^{\pi^{t+1}}(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$



Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^t} - V^{\star}\|_{\infty}$

$$\begin{aligned} V^{\star}(s) - V^{\pi^{t+1}}(s) &= \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \end{aligned}$$

$\pi^{t+1}(s)$
 $= \arg \max_a Q^{\pi^t}(s, a)$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

$$\begin{aligned} V^*(s) - V^{\pi^{t+1}}(s) &= \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \\ &= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^*(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s')) \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

$$\begin{aligned} V^*(s) - V^{\pi^{t+1}}(s) &= \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \\ &= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^*(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s')) \\ &\leq \max_a \left(\cancel{r(s, a)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')} - \left(\cancel{r(s, a)} + \underbrace{\gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s')} \right) \right) \end{aligned}$$

Analysis of Policy Iteration

Recall: Policy Improvement $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$ $\geq \gamma \|V^{\pi^{t+1}} - V^{\pi^t}\|_\infty$

Theorem: Convergence $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

$$V^*(s) - V^{\pi^{t+1}}(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right]$$

$$\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right]$$

$$= \max_a (r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^*(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s'))$$

$$\leq \max_a \left(\underbrace{r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')} - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \right)$$

$$\leq \gamma \|V^* - V^{\pi^t}\|_\infty \quad \leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} (V^*(s') - V^{\pi^t}(s'))$$

$$\leq \max_a \max_{s'} \gamma (V^*(s') - V^{\pi^t}(s'))$$

Analysis of Policy Iteration

Q: what happens when π^{t+1} and π^t are exactly the same?

Show that π^t is an optimal policy π^*

$$V^{\pi^t}(s) = V^*(s),$$

WS

Q: does this imply that the algorithm will terminate?

A S

Outline

1. Policy Iteration

2. Computation complexity of VI and PI

3. Linear Programming formulation

Computation complexity of VI and PI

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we **exactly** compute Q^* (or find π^*)
in time polynomial wrt $S, A, 1/(1 - \gamma)$?

Computation complexity of VI and PI

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we **exactly** compute Q^* (or find π^*)
in time polynomial wrt $S, A, 1/(1 - \gamma)$?

No for VI (i.e., gap between second and best)

Computation complexity of VI and PI

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we **exactly** compute Q^* (or find π^*)
in time polynomial wrt $S, A, 1/(1 - \gamma)$?

No for VI (i.e., gap between second and best)

Yes for policy iteration:

$$(S^3 + S^2A) \cdot \min \left\{ \frac{A^S}{S}, \frac{S^2A \log \frac{S^2}{1-\gamma}}{1-\gamma} \right\}$$

Complexity per iteration

Computation complexity of VI and PI

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$ can we **exactly** compute Q^\star (or find π^\star)
in time polynomial wrt $S, A, 1/(1 - \gamma)$?

No for VI (i.e., gap between second and best)

Yes for policy iteration:

$$(S^3 + S^2A) \cdot \min \left\{ \frac{A^S}{S}, \frac{S^2A \log \frac{S^2}{1-\gamma}}{1-\gamma} \right\}$$

What about $\text{poly}(S, A)$ algs?

Outline

1. Policy Iteration
2. Computation complexity of VI and PI
3. Linear Programming formulation

$$A^{\pi} \\ = Q^{\pi}(s,a) - V^{\pi}(s)$$

$$\max_a A^{\pi}(s,a)$$

The primal linear programming

Recall the Bellman consistency:

$$V(s) = \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V(s')] \right\}, \forall s$$

$$\Rightarrow V(s) = V^*(s), \forall s$$

$$V(s) \geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s'), \forall a$$

The primal linear programming

Recall the Bellman consistency:

$$V(s) = \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V(s')] \right\}, \forall s$$

We can re-write this as a linear program

The primal linear programming

Recall the Bellman consistency:

$$V(s) = \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V(s')] \right\}, \forall s$$

We can re-write this as a linear program

$$\begin{aligned} \min \sum_s \mu(s) V(s) & \quad \mu(s) > 0 \\ \text{s.t. } V(s) & \geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \quad \forall s, a \in S \times A \end{aligned}$$

$$V(s) > r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s'), \forall s$$

The primal linear programming

Recall the Bellman consistency:

$$V(s) = \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V(s')] \right\}, \forall s$$

We can re-write this as a linear program

$$\min \sum_s \mu(s) V(s)$$

$$\text{s.t. } V(s) \geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \quad \forall s, a \in S \times A$$

(Proof in HW1)

LP Runtime

[Ye, '05]: there is an interior point algorithm (CIPA)
which is (“nearly”) **strongly polynomial, i.e., no poly dependence on $1/(1 - \gamma)$**

$$S^4 A^4 \ln \left(\frac{S}{1 - \gamma} \right)$$

What about the Dual LP?

What about the Dual LP?

- Let us now consider the dual LP.
 - It is also very helpful conceptually.
 - In some cases, it also provides a reasonable algorithmic approach

What about the Dual LP?

- Let us now consider the dual LP.
 - It is also very helpful conceptually.
 - In some cases, it also provides a reasonable algorithmic approach
- Let us start by understanding the dual variables

State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of π visiting (s, a) at time step $h \in \mathbb{N}$, starting at s_0



State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of π visiting (s, a) at time step $h \in \mathbb{N}$, starting at s_0

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h(s, a; s_0, \pi)$$

$$\sum_{s,a} d_{s_0}^\pi(s, a) = 1$$



$h=0$

γ^0



$h=1$

γ^1

...



$h=\infty$

γ^∞

State action occupancy measure

$\mathbb{P}_h(s, a; s_0, \pi)$: probability of π visiting (s, a) at time step $h \in \mathbb{N}$, starting at s_0

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h(s, a; s_0, \pi)$$

$$V^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^\pi(s, a) r(s, a)$$

$$s_0 \sim \mu$$

$$\frac{1}{1 - \gamma} \mathbb{E}_{s_0 \sim d_{s_0}^\pi} [r(s, a)]$$

$$d_{\mu}^\pi(s, a) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s, a)]$$

A Bellman equation like property for $d_{s_0}^\pi(s, a)$

$$\sum_a d_\mu^\pi(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{\bar{s}, \bar{a}} P(s | \bar{s}, \bar{a}) d_\mu^\pi(\bar{s}, \bar{a})$$

Proof:

$$d_\mu^\pi(s, a) = (1 - \gamma) \left[\mu(s)\pi(a|s) + \gamma \boxed{P_1(s, a; \mu, \pi)} + \gamma^2 P_2(s, a; \mu, \pi) + \dots \right]$$

$$= (1 - \gamma) \left[\mu(s)\pi(a|s) + \gamma \underbrace{E_{\bar{s}, \bar{a} \sim P_0(\cdot; \mu, \pi)} [P(s | \bar{s}, \bar{a})\pi(a|s)]}_{\text{circled}} \right]$$

$$= (1 - \gamma) \left[\mu(s)\pi(a|s) + \left(\gamma \underbrace{E_{\bar{s}, \bar{a} \sim P_0(\cdot; \mu, \pi)} [P(s | \bar{s}, \bar{a})\pi(a|s)]}_{\text{circled}} + \gamma^2 \underbrace{E_{\bar{s}, \bar{a} \sim P_1(\cdot; \mu, \pi)} [P(s | \bar{s}, \bar{a})\pi(a|s)]}_{\text{circled}} + \gamma^3 \underbrace{E_{\bar{s}, \bar{a} \sim P_2(\cdot; \mu, \pi)} [P(s | \bar{s}, \bar{a})\pi(a|s)]}_{\text{circled}} + \dots \right) \right]$$

The “State-Action” Polytope

- Let us define the **state-action polytope** K as follows:

$$K_\mu := \left\{ \begin{array}{l} (d \mid d \geq 0 \text{ and } d \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|} \text{ and } d(s, g) = 0 \\ \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s \mid s', a') d(s', a') \end{array} \right\}$$

The “State-Action” Polytope

- Let us define the **state-action polytope** K as follows:

$$K_\mu := \left\{ d \mid d \geq 0 \text{ and } \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s \mid s', a') d(s', a') \right\}$$

- This set precisely characterizes all state-action visitation distributions:

The “State-Action” Polytope

- Let us define the **state-action polytope** K as follows:

$$K_\mu := \left\{ d \mid d \geq 0 \text{ and } \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s \mid s', a') d(s', a') \right\}$$

- This set precisely characterizes all state-action visitation distributions:
Lemma: $d \in K_\mu$ if and only if there exists a (possibly randomized) policy π s.t. $d_\mu^\pi = d$

The “State-Action” Polytope

- Let us define the **state-action polytope** K as follows:

$$K_\mu := \left\{ d \mid d \geq 0 \text{ and } \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s \mid s', a') d(s', a') \right\}$$

- This set precisely characterizes all state-action visitation distributions:
Lemma: $d \in K_\mu$ if and only if there exists a (possibly randomized) policy π

$$\text{s.t. } d_\mu^\pi = d$$

(Proof in HW1)

The Dual LP

$$\begin{aligned} \max_d \quad & \sum_{s,a} d(s,a)r(s,a) \\ \text{s.t.} \quad & d \in K_\mu \end{aligned}$$

Handwritten notes: A red arrow points from the d in the constraint to the d in the objective. To the right, the handwritten text $d^T r$ is written above μ .

- One can verify that this is the dual of the primal LP.

Summary

Notations: Value / Q functions, state-action occupant measures,
Bellman equation / optimality

Planning algorithms: VI, PI, LP (primal and dual)