

Reversible MDP's : Potential Applications?

Alejandro Gomez-Leos

UT Austin, Analysis and Design of Comm. Networks

December 5, 2023

Overview

- 1 Formulation [Barto+Sutton, 2018], [Puterman, 2014]
 - An example MDP
 - Basic definitions
 - Avg. reward criterion/policy gain

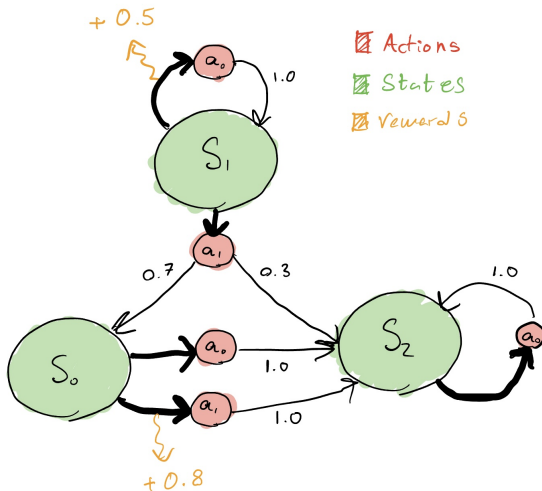
- 1 Formulation [Barto+Sutton, 2018], [Puterman, 2014]
 - An example MDP
 - Basic definitions
 - Avg. reward criterion/policy gain
- 2 Irreducible/reversible MDPs [Anantharam 2022]

- ① Formulation [Barto+Sutton, 2018], [Puterman, 2014]
 - An example MDP
 - Basic definitions
 - Avg. reward criterion/policy gain
- ② Irreducible/reversible MDPs [Anantharam 2022]
- ③ LP formulation/Poisson's equation, [Ross 1983], [Cogill+Peng, 2013]

- ① Formulation [Barto+Sutton, 2018], [Puterman, 2014]
 - An example MDP
 - Basic definitions
 - Avg. reward criterion/policy gain
- ② Irreducible/reversible MDPs [Anantharam 2022]
- ③ LP formulation/Poisson's equation, [Ross 1983], [Cogill+Peng, 2013]
- ④ Policy iteration [Anantharam 2022], [Cogill+Peng, 2013]

- ① Formulation [Barto+Sutton, 2018], [Puterman, 2014]
 - An example MDP
 - Basic definitions
 - Avg. reward criterion/policy gain
- ② Irreducible/reversible MDPs [Anantharam 2022]
- ③ LP formulation/Poisson's equation, [Ross 1983], [Cogill+Peng, 2013]
- ④ Policy iteration [Anantharam 2022], [Cogill+Peng, 2013]
- ⑤ Future Directions

An Example



Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

- \mathcal{S} for states, \mathcal{A} for actions

Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

- \mathcal{S} for states, \mathcal{A} for actions
- $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}|}$ is state-action-state transition matrix

Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

- \mathcal{S} for states, \mathcal{A} for actions
- $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ is state-action-state transition matrix
- $\mathbf{r} \in [0, 1]^{|\mathcal{S}||\mathcal{A}|}$ is reward vector

Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

- \mathcal{S} for states, \mathcal{A} for actions
- $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ is state-action-state transition matrix
- $\mathbf{r} \in [0, 1]^{|\mathcal{S}||\mathcal{A}|}$ is reward vector

① Objects collectively called environment

Definition

A MDP is a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$

- \mathcal{S} for states, \mathcal{A} for actions
- $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ is state-action-state transition matrix
- $\mathbf{r} \in [0, 1]^{|\mathcal{S}||\mathcal{A}|}$ is reward vector

- 1 Objects collectively called environment
- 2 Stochastic process arises by fixing a **policy**

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

- Π_{det} : set of deterministic policies for \mathcal{M}

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

- Π_{det} : set of deterministic policies for \mathcal{M}
- Π_{rand} : set of randomized policies \mathcal{M}

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

- Π_{det} : set of deterministic policies for \mathcal{M}
- Π_{rand} : set of randomized policies \mathcal{M}
- $\Pi_{\text{det}} \subseteq \Pi_{\text{rand}}$

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

- Π_{det} : set of deterministic policies for \mathcal{M}
- Π_{rand} : set of randomized policies \mathcal{M}
- $\Pi_{\text{det}} \subseteq \Pi_{\text{rand}}$
- $\pi(a \mid s)$: prob. of taking $a \in \mathcal{A}$ given in state s

Definition

A policy $\pi : \mathcal{S} \rightarrow \Delta^{|\mathcal{A}|}$ is a conditional distribution over \mathcal{A} for each $s \in \mathcal{S}$.

Notation:

- Π_{det} : set of deterministic policies for \mathcal{M}
- Π_{rand} : set of randomized policies \mathcal{M}
- $\Pi_{\text{det}} \subseteq \Pi_{\text{rand}}$
- $\pi(a \mid s)$: prob. of taking $a \in \mathcal{A}$ given in state s
- If $\pi \in \Pi_{\text{det}}$, $\pi(s)$ is action for state s

Objective

Objective: The agent's goal is to pick a policy π , such that it "accrues maximal reward."

Objective

Objective: The agent's goal is to pick a policy π , such that it "accrues maximal reward."

Several notions of "maximal reward"—different objectives

Theorem (Puterman, 2005)

A policy $\pi \in \Pi_{rand}$ is associated with a canonical measure on the sequence of r.v.'s $S_0, A_0, S_1, A_1, \dots$ for the MDP \mathcal{M} .

Theorem (Puterman, 2005)

A policy $\pi \in \Pi_{rand}$ is associated with a canonical measure on the sequence of r.v.'s $S_0, A_0, S_1, A_1, \dots$ for the MDP \mathcal{M} .

The performance measure:

- $v_{avg}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right)$, for $\mathbf{u} \in \Delta^{|S|}$ aka gain

Theorem (Puterman, 2005)

A policy $\pi \in \Pi_{rand}$ is associated with a canonical measure on the sequence of r.v.'s $S_0, A_0, S_1, A_1, \dots$ for the MDP \mathcal{M} .

The performance measure:

- $v_{avg}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right)$, for $\mathbf{u} \in \Delta^{|S|}$ aka gain

What is a reasonable objective?

Theorem (Puterman, 2005)

A policy $\pi \in \Pi_{rand}$ is associated with a canonical measure on the sequence of r.v.'s $S_0, A_0, S_1, A_1, \dots$ for the MDP \mathcal{M} .

The performance measure:

- $v_{avg}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right)$, for $\mathbf{u} \in \Delta^{|S|}$ aka gain

What is a reasonable objective?

- Find π to max $v_{avg}^{\pi}(\mathbf{u})$ for all \mathbf{u} ?

Theorem (Puterman, 2005)

A policy $\pi \in \Pi_{rand}$ is associated with a canonical measure on the sequence of r.v.'s $S_0, A_0, S_1, A_1, \dots$ for the MDP \mathcal{M} .

The performance measure:

- $v_{avg}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} r_{S_t, A_t} \right)$, for $\mathbf{u} \in \Delta^{|S|}$ aka gain

What is a reasonable objective?

- Find π to max $v_{avg}^{\pi}(\mathbf{u})$ for all \mathbf{u} ?
- Maybe study $v_{avg}^{\pi}(\mathbf{u}) = v_{avg}^{\pi}(\mathbf{u}') \implies$ notion of irreducibility needed.

Irreducible and reversible MDPs

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Definition (Irreducible MDP)

If $\forall \pi \in \Pi_{\text{rand}}$, $(\mathcal{S}, \mathbf{P}^\pi)$ irreducible, then \mathcal{M} is irreducible.

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Definition (Irreducible MDP)

If $\forall \pi \in \Pi_{\text{rand}}$, $(\mathcal{S}, \mathbf{P}^\pi)$ irreducible, then \mathcal{M} is irreducible.

Example:

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Definition (Irreducible MDP)

If $\forall \pi \in \Pi_{\text{rand}}, (\mathcal{S}, \mathbf{P}^\pi)$ irreducible, then \mathcal{M} is irreducible.

Example: Any MDP with $\mathbf{P}_{(s,a),s'} > 0$ for all (s, a, s')

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Definition (Irreducible MDP)

If $\forall \pi \in \Pi_{\text{rand}}, (\mathcal{S}, \mathbf{P}^\pi)$ irreducible, then \mathcal{M} is irreducible.

Example: Any MDP with $\mathbf{P}_{(s,a),s'} > 0$ for all (s, a, s')

Definition (Reversible MDP, Cogill+Peng, 2013, Anantharam 2022)

If $\forall \pi \in \Pi_{\text{rand}}, (\mathcal{S}, \mathbf{P}^\pi)$ irreducible and reversible, then \mathcal{M} is reversible.

Irreducible and reversible MDPs

- The **on-policy** transition matrix \mathbf{P}^π for chain $(\mathcal{S}, \mathbf{P}^\pi)$ given by
$$\mathbf{P}_{s,s'}^\pi := \sum_a \mathbf{P}_{(s,a),s'} \pi(a | s)$$

Definition (Irreducible MDP)

If $\forall \pi \in \Pi_{\text{rand}}$, $(\mathcal{S}, \mathbf{P}^\pi)$ irreducible, then \mathcal{M} is irreducible.

Example: Any MDP with $\mathbf{P}_{(s,a),s'} > 0$ for all (s, a, s')

Definition (Reversible MDP, Cogill+Peng, 2013, Anantharam 2022)

If $\forall \pi \in \Pi_{\text{rand}}$, $(\mathcal{S}, \mathbf{P}^\pi)$ irreducible and reversible, then \mathcal{M} is reversible.

Note: Could replace Π_{rand} by Π_{det} in def's above.

Reversible MDP Example

Theorem (Anantharam, 2022)

Reversible MDP Example

Theorem (Anantharam, 2022)

- 1 *Let \mathcal{S}, \mathcal{A} be any finite sets*

Reversible MDP Example

Theorem (Anantharam, 2022)

- 1 Let S, \mathcal{A} be any finite sets
- 2 Let $\mathcal{G} = (S, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let S, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (S, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $Q_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let \mathcal{S}, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $Q_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- ④ Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let \mathcal{S}, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $\mathbf{Q}_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- ④ Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$
- ⑤ $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a) \mathbf{Q}_{s,s'} \quad \mathbf{r}_{s,a} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let \mathcal{S}, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $\mathbf{Q}_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- ④ Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$
- ⑤ $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a) \mathbf{Q}_{s,s'} \quad \mathbf{r}_{s,a} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$

$\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$ is a reversible MDP.

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let \mathcal{S}, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $\mathbf{Q}_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- ④ Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$
- ⑤ $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a) \mathbf{Q}_{s,s'} \quad \mathbf{r}_{s,a} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$

$\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$ is a reversible MDP.

Interpretation:

Reversible MDP Example

Theorem (Anantharam, 2022)

- 1 Let \mathcal{S}, \mathcal{A} be any finite sets
- 2 Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- 3 $\mathbf{Q}_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- 4 Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$
- 5 $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a) \mathbf{Q}_{s,s'} \quad \mathbf{r}_{s,a} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$

$\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$ is a reversible MDP.

Interpretation:

- like *lazy weighted random walk*, with variable laziness via ρ

Reversible MDP Example

Theorem (Anantharam, 2022)

- ① Let \mathcal{S}, \mathcal{A} be any finite sets
- ② Let $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ any simple connected graph, with
$$(s, s') \in \mathcal{E} \iff \exists w_{s,s'} = w_{s',s} > 0$$
- ③ $\mathbf{Q}_{s,s'} = \frac{w_{s,s'}}{w_s}$ for all (s, s') , where $w_s := \sum_{s'} w_{s,s'}$
- ④ Let $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$
- ⑤ $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a) \mathbf{Q}_{s,s'} \quad \mathbf{r}_{s,a} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$

$\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r})$ is a reversible MDP.

Interpretation:

- like *lazy weighted random walk*, with variable laziness via ρ
- actions can control laziness, but not totally

A more general characterization

- This construction exhibits a reversible MDP

A more general characterization

- This construction exhibits a reversible MDP
- Conversely, do all reversible MDPs have that structure?

A more general characterization

- This construction exhibits a reversible MDP
- Conversely, do all reversible MDPs have that structure? Sort of!

A more general characterization

- This construction exhibits a reversible MDP
- Conversely, do all reversible MDPs have that structure? Sort of!

Lemma (Anatharam, '22)

Consider reversible MDP \mathcal{M} . There exists a simple connected graph $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ s.t. for all $a \in \mathcal{A}$ and $s \neq s'$, $P_{(s,a),s'} > 0$ iff $(s, s') \in E$.

A more general characterization

- This construction exhibits a reversible MDP
- Conversely, do all reversible MDPs have that structure? Sort of!

Lemma (Anatharam, '22)

Consider reversible MDP \mathcal{M} . There exists a simple connected graph $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ s.t. for all $a \in \mathcal{A}$ and $s \neq s'$, $P_{(s,a),s'} > 0$ iff $(s, s') \in E$.

Theorem ("bi-connection theorem" Anatharam, '22)

Let \mathcal{M} be a reversible MDP, and $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ its canonical graph. If \mathcal{G} is *bi-connected*, then there exists a irr. and reversible \mathbf{Q} and a function $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ such that for each $a \in \mathcal{A}$,

- $\mathbf{Q}_{s,s'} > 0$ iff $(s, s') \in \mathcal{E}$
- $\mathbf{P}_{(s,a),s} := 1 - \rho(s, a), \quad \mathbf{P}_{(s,a),s'} := \rho(s, a)\mathbf{Q}_{s,s'}$

Optimization

Back to optimization.

Back to optimization.

Suppose \mathcal{M} is irreducible, recall:

Back to optimization.

Suppose \mathcal{M} is irreducible, recall:

$$v_{\text{avg}}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right), \text{ for } \mathbf{u} \in \Delta^{|S|}$$

Back to optimization.

Suppose \mathcal{M} is irreducible, recall:

$$v_{\text{avg}}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right), \text{ for } \mathbf{u} \in \Delta^{|\mathcal{S}|}$$

For any π , $(\mathcal{S}, \mathbf{P}^{\pi})$ is irreducible, so sum converges to same quantity regardless of \mathbf{u} . One can show

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \mu^{\pi}(s) \pi(a | s) \mathbf{r}_{s,a}$$

Back to optimization.

Suppose \mathcal{M} is irreducible, recall:

$$v_{\text{avg}}^{\pi}(\mathbf{u}) := \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi} \left(\sum_{t=0}^{T-1} \mathbf{r}_{S_t, A_t} \right), \text{ for } \mathbf{u} \in \Delta^{|\mathcal{S}|}$$

For any π , $(\mathcal{S}, \mathbf{P}^{\pi})$ is irreducible, so sum converges to same quantity regardless of \mathbf{u} . One can show

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \mu^{\pi}(s) \pi(a | s) \mathbf{r}_{s,a}$$

How to pick π to maximize?

Primal Formulation

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \underbrace{\mu^{\pi}(s) \pi(a \mid s)}_{:= \phi^{\pi}(s,a)} \mathbf{r}_{s,a}$$

Primal Formulation

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \underbrace{\mu^{\pi}(s) \pi(a \mid s)}_{:= \phi^{\pi}(s,a)} \mathbf{r}_{s,a}$$

Solve the following LP OPT:

$$\begin{aligned} \max_{\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \quad & \sum_{s,a} \phi(s, a) \mathbf{r}_{s,a} \\ \text{s.t.} \quad & \phi(s, a) \geq 0 \quad \forall s, a \\ & \sum_{s,a} \phi(s, a) = 1 \\ & \sum_{s,a} \phi(s, a) = \sum_{s',a'} \mathbf{P}_{(s',a'),s} \phi(s', a') \quad \forall s \end{aligned}$$

Primal Formulation

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \underbrace{\mu^{\pi}(s) \pi(a | s)}_{:= \phi^{\pi}(s,a)} \mathbf{r}_{s,a}$$

Solve the following LP OPT:

$$\begin{aligned} \max_{\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \quad & \sum_{s,a} \phi(s, a) \mathbf{r}_{s,a} \\ \text{s.t.} \quad & \phi(s, a) \geq 0 \quad \forall s, a \\ & \sum_{s,a} \phi(s, a) = 1 \\ & \sum_{s,a} \phi(s, a) = \sum_{s',a'} \mathbf{P}_{(s',a'),s} \phi(s', a') \quad \forall s \end{aligned}$$

Note: feasible thanks to irreducibility of \mathcal{M} —but, does there exists an optimal solution?

Primal Formulation

$$v_{\text{avg}}^{\pi} = \sum_{s,a} \underbrace{\mu^{\pi}(s) \pi(a | s)}_{:= \phi^{\pi}(s,a)} \mathbf{r}_{s,a}$$

Solve the following LP OPT:

$$\begin{aligned} \max_{\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \quad & \sum_{s,a} \phi(s, a) \mathbf{r}_{s,a} \\ \text{s.t.} \quad & \phi(s, a) \geq 0 \quad \forall s, a \\ & \sum_{s,a} \phi(s, a) = 1 \\ & \sum_{s,a} \phi(s, a) = \sum_{s',a'} \mathbf{P}_{(s',a'),s} \phi(s', a') \quad \forall s \end{aligned}$$

Note: feasible thanks to irreducibility of \mathcal{M} —but, does there exists an optimal solution? **Yes!**

Dual LP of OPT

$$\begin{aligned} \max_{\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}} \quad & \sum_{s,a} \phi(s,a) r_{s,a} \\ \text{s.t.} \quad & \phi(s,a) \geq 0 \quad \forall s,a \\ & \sum_{s,a} \phi(s,a) = 1 \\ & \sum_{s,a} \phi(s,a) = \sum_{s',a'} P_{(s',a'),s} \phi(s',a') \quad \forall s \end{aligned}$$

OPT admits the dual LP:

$$\begin{aligned} \min_{h: \mathcal{S} \rightarrow \mathbb{R}, C \in \mathbb{R}} \quad & C \\ \text{s.t.} \quad & r_{s,a} + \sum_{s'} h(s') P_{(s,a),s'} \leq h(s) + C \quad \forall s,a \end{aligned}$$

Existence of optimal det. policy

Claim: We can find a $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT.

Existence of optimal det. policy

Claim: We can find $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT.

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that

$$C + h(s) = \max_a \left(r_{s,a} + \sum_{s'} \mathbf{P}_{(s,a),s'} h(s') \right) \quad \forall s \in \mathcal{S}$$

then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Existence of optimal det. policy

Claim: We can find $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT.

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that

$$C + h(s) = \max_a \left(r_{s,a} + \sum_{s'} \mathbf{P}_{(s,a),s'} h(s') \right) \quad \forall s \in \mathcal{S}$$

then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Notation: For more compact formulas, define

$$\psi^h(s, a) := r_{s,a} + \sum_{s'} \mathbf{P}_{(s,a),s'} h(s')$$

for each given h and pair (s, a)

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Continuing

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Pick (h^*, C^*) as dual optimal sol'ns. Goal is to exhibit they satisfy lemma.

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Pick (h^*, C^*) as dual optimal sol'ns. Goal is to exhibit they satisfy lemma.

For each s , let $a_s := \arg \max_a \psi^{h^*}(s, a)$.

Continuing

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Pick (h^*, C^*) as dual optimal sol'ns. Goal is to exhibit they satisfy lemma.

For each s , let $a_s := \arg \max_a \psi^{h^*}(s, a)$.

BWOC there exists \tilde{s} s.t. $\psi^{h^*}(\tilde{s}, a_{\tilde{s}}) < h^*(\tilde{s}) + C^*$.

Continuing

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Pick (h^*, C^*) as dual optimal sol'ns. Goal is to exhibit they satisfy lemma.

For each s , let $a_s := \arg \max_a \psi^{h^*}(s, a)$.

BWOC there exists \tilde{s} s.t. $\psi^{h^*}(\tilde{s}, a_{\tilde{s}}) < h^*(\tilde{s}) + C^*$.

By C.S. $\phi^*(\tilde{s}, a_{\tilde{s}}) = 0$, so $\exists a' : \phi^*(\tilde{s}, a') > 0 \implies \psi^{h^*}(\tilde{s}, a') = h^*(\tilde{s}) + C^*$.

Continuing

Lemma (S.M. Ross, 1983)

For any irreducible MDP \mathcal{M} , if there is a bounded real function $h : \mathcal{S} \rightarrow \mathbb{R}$, and a constant C , such that $C + h(s) = \max_a \psi^h(s, a)$ for all $s \in \mathcal{S}$, then there exists an optimal $\pi \in \Pi_{\text{det}}$, where $v_{\text{avg}}^\pi = C$.

Claim: There exists a policy $\pi \in \Pi_{\text{det}}$ such that $\phi^\pi = \phi^*$ for OPT

Proof.

Pick (h^*, C^*) as dual optimal sol'ns. Goal is to exhibit they satisfy lemma.

For each s , let $a_s := \arg \max_a \psi^{h^*}(s, a)$.

BWOC there exists \tilde{s} s.t. $\psi^{h^*}(\tilde{s}, a_{\tilde{s}}) < h^*(\tilde{s}) + C^*$.

By C.S. $\phi^*(\tilde{s}, a_{\tilde{s}}) = 0$, so $\exists a' : \phi^*(\tilde{s}, a') > 0 \implies \psi^{h^*}(\tilde{s}, a') = h^*(\tilde{s}) + C^*$.

But $h^*(\tilde{s}) + C^* = \psi^{h^*}(\tilde{s}, a') \leq \psi^{h^*}(\tilde{s}, a_{\tilde{s}}) < h^*(\tilde{s}) + C^*$, contradiction.

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- 2 (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal).

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- 2 (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal). Fix an $s \in \mathcal{S}$.

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- 2 (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal). Fix an $s \in \mathcal{S}$. The following holds iff (h^*, C^*) are dual-optimal by S.D.

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- 2 (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal). Fix an $s \in \mathcal{S}$. The following holds iff (h^*, C^*) are dual-optimal by S.D.

- 1 If $a \neq \pi^*(s)$, then $\phi^*(s, a) = 0 \implies \psi^{h^*}(s, a) < h^*(s) + C^*$

Dual LP of OPT properties

Facts:

- 1 $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- 2 (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal). Fix an $s \in \mathcal{S}$. The following holds iff (h^*, C^*) are dual-optimal by S.D.

- 1 If $a \neq \pi^*(s)$, then $\phi^*(s, a) = 0 \implies \psi^{h^*}(s, a) < h^*(s) + C^*$
- 2 If $a = \pi^*(s)$, then $\phi^*(s, a) > 0 \implies \psi^{h^*}(s, a) = h^*(s) + C^*$

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Proof.

Take $\phi^* = \phi^{\pi^*}$ (this is OK because ϕ^{π^*} is primal-optimal). Fix an $s \in \mathcal{S}$. The following holds iff (h^*, C^*) are dual-optimal by S.D.

- ① If $a \neq \pi^*(s)$, then $\phi^*(s, a) = 0 \implies \psi^{h^*}(s, a) < h^*(s) + C^*$
- ② If $a = \pi^*(s)$, then $\phi^*(s, a) > 0 \implies \psi^{h^*}(s, a) = h^*(s) + C^*$

Hence $\max_a \psi^{h^*}(s, a) = h^*(s) + C^*$.

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration:

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration: Input initial $\pi_0 \in \Pi_{\text{det}}$, for $k = 0, 1, 2, \dots$

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration: Input initial $\pi_0 \in \Pi_{\text{det}}$, for $k = 0, 1, 2, \dots$

- Solve *Poisson's eq'n*

$$\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

for (h^{π_k}, C^{π_k}) (system of $|\mathcal{S}|$ eq'ns in $|\mathcal{S}| + 1$ unknowns)

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration: Input initial $\pi_0 \in \Pi_{\text{det}}$, for $k = 0, 1, 2, \dots$

- Solve *Poisson's eq'n*

$$\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

for (h^{π_k}, C^{π_k}) (system of $|\mathcal{S}|$ eq'ns in $|\mathcal{S}| + 1$ unknowns)

- Check if

$$\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration: Input initial $\pi_0 \in \Pi_{\text{det}}$, for $k = 0, 1, 2, \dots$

- Solve *Poisson's eq'n*

$$\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

for (h^{π_k}, C^{π_k}) (system of $|\mathcal{S}|$ eq'ns in $|\mathcal{S}| + 1$ unknowns)

- Check if

$$\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

- If not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Dual LP of OPT properties

Facts:

- ① $\exists \pi^* \in \Pi_{\text{det}}$ given by $\pi^*(s) = \arg \max_a \psi^{h^*}(s, a)$, for which $\phi^{\pi^*} = \phi^*$
- ② (h^*, C^*) are dual-optimal iff $\max_a \psi^{h^*}(s, a) = h^*(s) + C$ for all s .

Policy Iteration: Input initial $\pi_0 \in \Pi_{\text{det}}$, for $k = 0, 1, 2, \dots$

- Solve *Poisson's eq'n*

$$\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

for (h^{π_k}, C^{π_k}) (system of $|\mathcal{S}|$ eq'ns in $|\mathcal{S}| + 1$ unknowns)

- Check if

$$\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$$

- If not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Next page: Proof that **update strictly improves objective**.

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$v_{\text{avg}}^{\pi_k} = C^{\pi_k}$$

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$v_{\text{avg}}^{\pi_k} = C^{\pi_k} = \sum_s \mu^{\pi_{k+1}}(s) C^{\pi_k}$$

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$\begin{aligned} v_{\text{avg}}^{\pi_k} &= C^{\pi_k} = \sum_s \mu^{\pi_{k+1}}(s) C^{\pi_k} \\ &< \sum_s \mu^{\pi_{k+1}}(s) [\psi^{\pi_k}(s, \pi_{k+1}(s)) - h^{\pi_k}(s)] \end{aligned}$$

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$\begin{aligned} v_{\text{avg}}^{\pi_k} &= C^{\pi_k} = \sum_s \mu^{\pi_{k+1}}(s) C^{\pi_k} \\ &< \sum_s \mu^{\pi_{k+1}}(s) [\psi^{\pi_k}(s, \pi_{k+1}(s)) - h^{\pi_k}(s)] \\ &= \sum_s \mu^{\pi_{k+1}}(s) [\mathbf{r}_{s, \pi_{k+1}(s)} + \sum_{s'} \mathbf{P}_{(s, \pi_{k+1}(s)), s'} h^{\pi_k}(s') - h^{\pi_k}(s)] \end{aligned}$$

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$\begin{aligned} v_{\text{avg}}^{\pi_k} &= C^{\pi_k} = \sum_s \mu^{\pi_{k+1}}(s) C^{\pi_k} \\ &< \sum_s \mu^{\pi_{k+1}}(s) [\psi^{\pi_k}(s, \pi_{k+1}(s)) - h^{\pi_k}(s)] \\ &= \sum_s \mu^{\pi_{k+1}}(s) [\mathbf{r}_{s, \pi_{k+1}(s)} + \sum_{s'} \mathbf{P}_{(s, \pi_{k+1}(s)), s'} h^{\pi_k}(s') - h^{\pi_k}(s)] \\ &= v_{\text{avg}}^{\pi_{k+1}} + \sum_s \left(\sum_{s'} \mu^{\pi_{k+1}}(s) \mathbf{P}_{(s, \pi_{k+1}(s)), s'} h^{\pi_k}(s') - \mu^{\pi_{k+1}}(s) h^{\pi_k}(s) \right) \end{aligned}$$

Monotonic Improvement

Policy Iteration: For a deterministic input policy π_k , first

- solve *Poisson's eq'n* $\psi^{h^{\pi_k}}(s, \pi_k(s)) = h^{\pi_k}(s) + C^{\pi_k}$ for (h^{π_k}, C^{π_k})
- check if $\max_a \psi^{h^{\pi_k}}(s, a) = h^{\pi_k}(s) + C^{\pi_k} \quad \forall s \in \mathcal{S}$
- if not, define updated policy $\pi_{k+1}(s) = \arg \max_a \psi^{h^{\pi_k}}(s, a)$

Claim: $v_{\text{avg}}^{\pi_k} < v_{\text{avg}}^{\pi_{k+1}}$, and procedure converges in finite steps.

$$\begin{aligned} v_{\text{avg}}^{\pi_k} &= C^{\pi_k} = \sum_s \mu^{\pi_{k+1}}(s) C^{\pi_k} \\ &< \sum_s \mu^{\pi_{k+1}}(s) [\psi^{\pi_k}(s, \pi_{k+1}(s)) - h^{\pi_k}(s)] \\ &= \sum_s \mu^{\pi_{k+1}}(s) [\mathbf{r}_{s, \pi_{k+1}(s)} + \sum_{s'} \mathbf{P}_{(s, \pi_{k+1}(s)), s'} h^{\pi_k}(s') - h^{\pi_k}(s)] \\ &= v_{\text{avg}}^{\pi_{k+1}} + \sum_s \left(\sum_{s'} \mu^{\pi_{k+1}}(s) \mathbf{P}_{(s, \pi_{k+1}(s)), s'} h^{\pi_k}(s') - \mu^{\pi_{k+1}}(s) h^{\pi_k}(s) \right) \\ &= v_{\text{avg}}^{\pi_{k+1}} \end{aligned}$$

Discussion

- Main computational burden is solving Poisson's equations.

- Main computational burden is solving Poisson's equations.
- We haven't used reversibility yet.

- Main computational burden is solving Poisson's equations.
- We haven't used reversibility yet.
- Reversibility leads to a simplified policy iteration.

The Theorem

Theorem (Anantharam, 2022)

The Theorem

Theorem (Anantharam, 2022)

*Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**.*

The Theorem

Theorem (Anantharam, 2022)

*Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed.*

The Theorem

Theorem (Anantharam, 2022)

*Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.*

The Theorem

Theorem (Anantharam, 2022)

*Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.*

Let some initial $\pi_0 \in \Pi_{det}$. For $k = 0, 1, 2 \dots$

The Theorem

Theorem (Anantharam, 2022)

*Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.*

Let some initial $\pi_0 \in \Pi_{det}$. For $k = 0, 1, 2 \dots$

- 1 *For π_k compute C^{π_k} .*

The Theorem

Theorem (Anantharam, 2022)

Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is *bi-connected*. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.

Let some initial $\pi_0 \in \Pi_{det}$. For $k = 0, 1, 2 \dots$

- 1 For π_k compute C^{π_k} .
- 2 If $\exists s$ such that

$$\frac{\mathbf{r}_{(s, \pi_k(s))} - C^{\pi_k}}{\rho(s, \pi_k(s))} < \arg \max_a \frac{\mathbf{r}_{(s, a)} - C^{\pi_k}}{\rho(s, a)} \quad (1)$$

Let a_s denote $\arg \max$.

The Theorem

Theorem (Anantharam, 2022)

Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is **bi-connected**. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.

Let some initial $\pi_0 \in \Pi_{det}$. For $k = 0, 1, 2 \dots$

- 1 For π_k compute C^{π_k} .
- 2 If $\exists s$ such that

$$\frac{\mathbf{r}_{(s, \pi_k(s))} - C^{\pi_k}}{\rho(s, \pi_k(s))} < \arg \max_a \frac{\mathbf{r}_{(s, a)} - C^{\pi_k}}{\rho(s, a)} \quad (1)$$

Let a_s denote $\arg \max$. Set $\pi_{k+1}(s) = a_s$, and $\pi_{k+1}(s') = \pi_k(s')$ for $s' \neq s$. If no such s , terminate.

The Theorem

Theorem (Anantharam, 2022)

Consider a reversible MDP \mathcal{M} , whose canonical graph \mathcal{G} is *bi-connected*. Let \mathbf{Q} and $\rho : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$ be as guaranteed. Consider the following iterative procedure.

Let some initial $\pi_0 \in \Pi_{det}$. For $k = 0, 1, 2 \dots$

- 1 For π_k compute C^{π_k} .
- 2 If $\exists s$ such that

$$\frac{\mathbf{r}_{(s, \pi_k(s))} - C^{\pi_k}}{\rho(s, \pi_k(s))} < \arg \max_a \frac{\mathbf{r}_{(s, a)} - C^{\pi_k}}{\rho(s, a)} \quad (1)$$

Let a_s denote $\arg \max$. Set $\pi_{k+1}(s) = a_s$, and $\pi_{k+1}(s') = \pi_k(s')$ for $s' \neq s$. If no such s , terminate.

Then $v_{avg}^{\pi_k} < v_{avg}^{\pi_{k+1}}$, and procedure converges in finite steps to an opt. policy.

Proof?

Main Takeaways

Main Takeaways

- Policy gain is an alternative objective function for policy opt.

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!
- **Future Directions:**

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!
- **Future Directions:**
 - Finite time bounds for ϵ -approximation of opt. policy? Lower bounds without reversibility? Upper bounds with?

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!
- **Future Directions:**
 - Finite time bounds for ϵ -approximation of opt. policy? Lower bounds without reversibility? Upper bounds with?
 - RL version—transition dynamics unknown prior, (ϵ, δ) guarantees...

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!
- **Future Directions:**
 - Finite time bounds for ϵ -approximation of opt. policy? Lower bounds without reversibility? Upper bounds with?
 - RL version—transition dynamics unknown prior, (ϵ, δ) guarantees...
Avg. reward RL still an active area of research!

Main Takeaways

- Policy gain is an alternative objective function for policy opt.
- Optimizing gain for MDPs seems quite tough
 - Dual formulation led to algorithm requiring $\mathcal{O}(|S|)$ linear system solver
- Under reversibility + a little more, could get a better algorithm
 - No linear system solving!
- **Future Directions:**
 - Finite time bounds for ϵ -approximation of opt. policy? Lower bounds without reversibility? Upper bounds with?
 - RL version—transition dynamics unknown prior, (ϵ, δ) guarantees...
Avg. reward RL still an active area of research!
 - Learnability—identification of irreducible or reversible MDP, learnability of ρ

Thanks!

- Ross, S. M. "Introduction to Stochastic Dynamic Programming, 1983."
- Cogill, Randy, and Cheng Peng. "Reversible Markov decision processes with an average-reward criterion." *SIAM Journal on Control and Optimization* 51.1 (2013): 402-418.
- Puterman, Martin L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley and Sons, 2014.
- Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Anantharam, Venkat. "Reversible Markov decision processes and the Gaussian free field." *Systems and Control Letters* 169 (2022): 105382.