

Provably Adaptive Average Reward Reinforcement Learning for Metric Spaces

Avik Kar
ECE

Rahul Singh
ECE



3rd April 2025

IISc, Bengaluru

This talk contains

Average Reward Reinforcement Learning

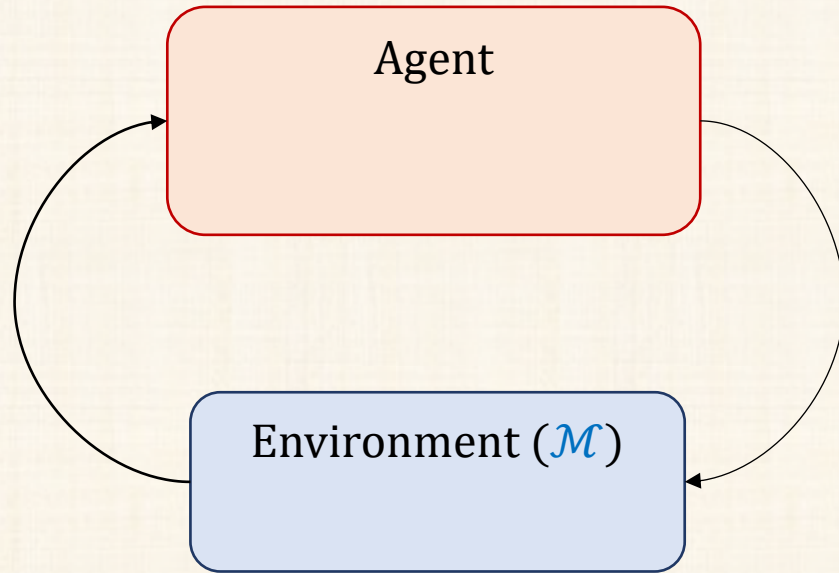
Lipschitz Continuity of Average Reward

Algorithm: Zooming in Policy Space

Algorithm: Zooming in State-Action Space

Reinforcement Learning Setup

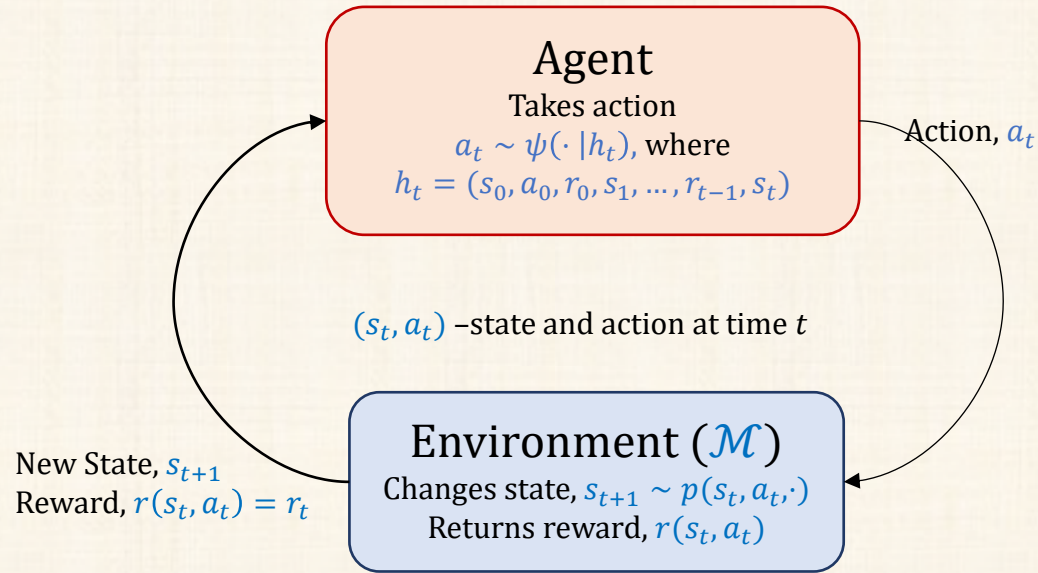
Reinforcement Learning Setup



Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

Reinforcement Learning Setup

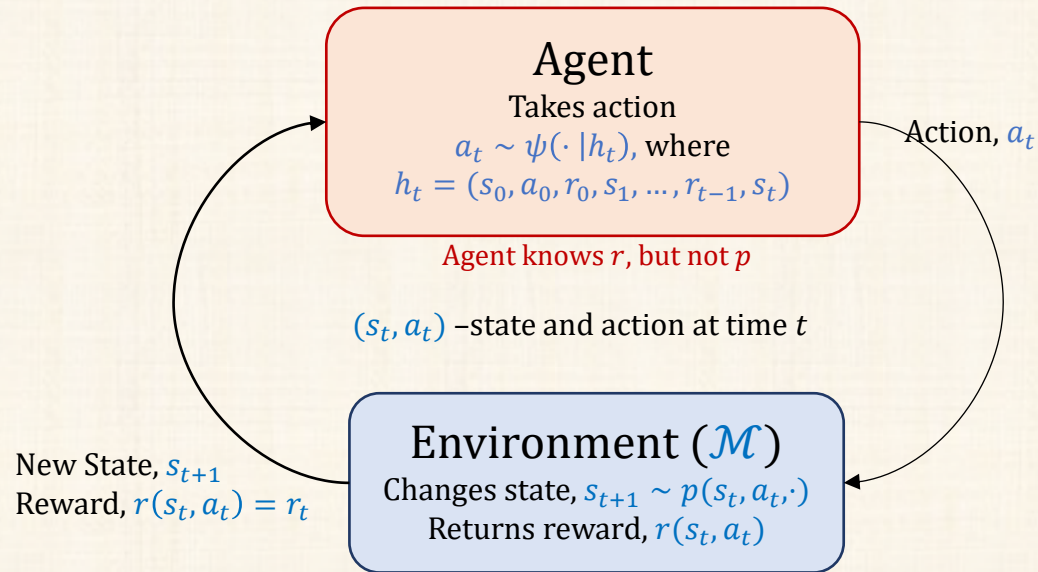


Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

Reinforcement Learning Setup

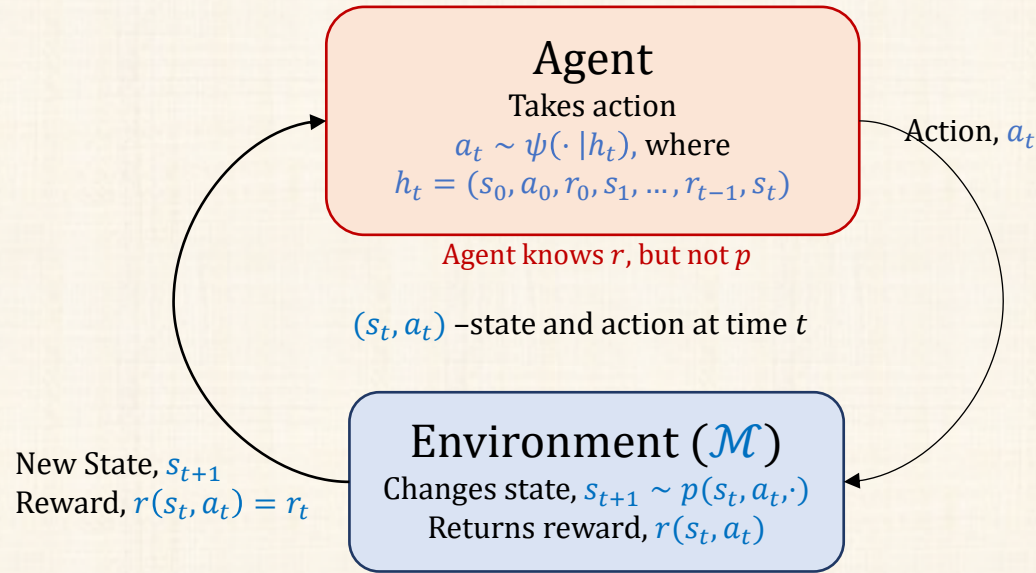


Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

Reinforcement Learning Setup



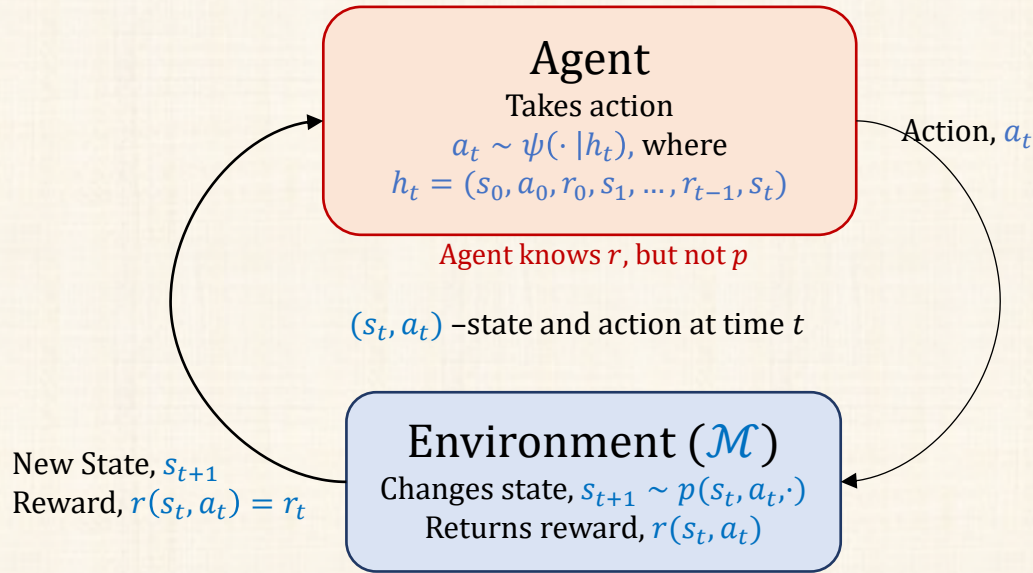
- Stationary deterministic policy: a function $\phi: \mathcal{S} \rightarrow \mathcal{A}$

Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

Reinforcement Learning Setup



- Stationary deterministic policy: a function $\phi: \mathcal{S} \rightarrow \mathcal{A}$

Average reward criterion:

- Let

$$J_{\mathcal{M}}(\phi) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\phi} \left[\sum_{t=1}^T r_t \right]$$

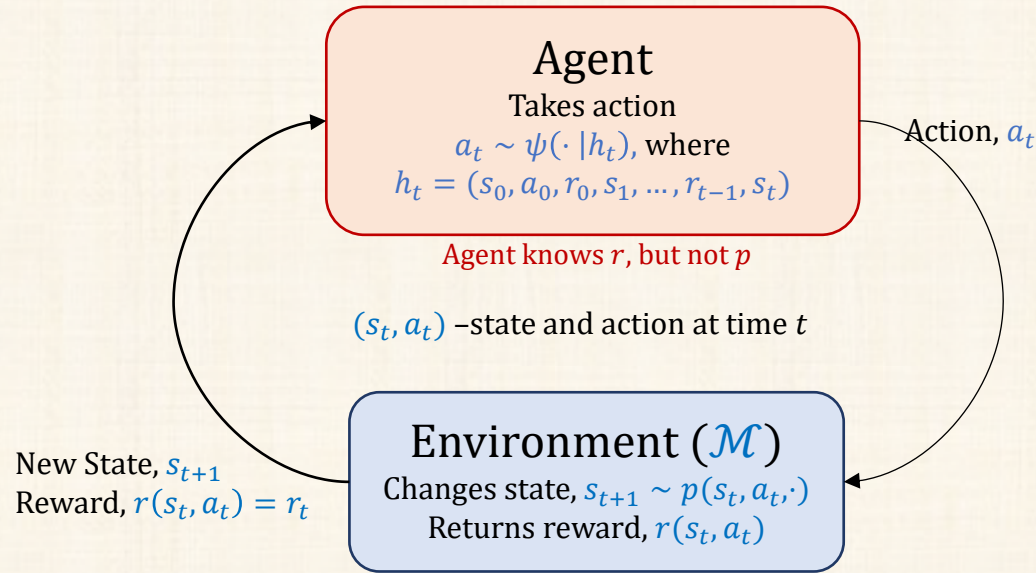
Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
 $p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

★ \mathbb{E}_{ϕ} – Expectation taken considering policy ϕ is played

Reinforcement Learning Setup



- Stationary deterministic policy: a function $\phi: \mathcal{S} \rightarrow \mathcal{A}$

Average reward criterion:

- Let

$$J_{\mathcal{M}}(\phi) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\phi} \left[\sum_{t=1}^T r_t \right]$$

- Optimal average reward

$$J_{\mathcal{M}}^* = \max_{\phi} J_{\mathcal{M}}(\phi)$$

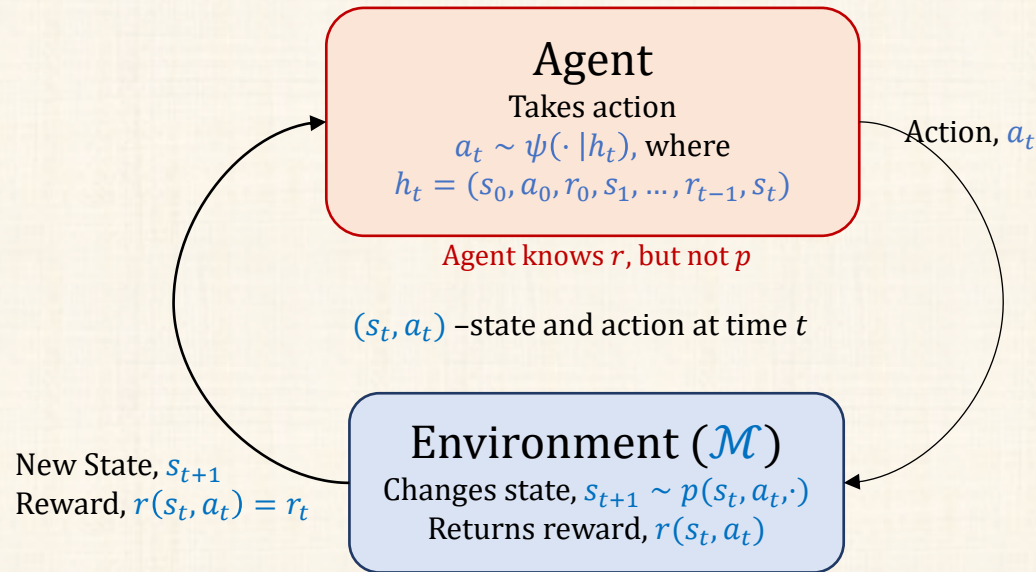
Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;
 $p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

★ \mathbb{E}_{ϕ} – Expectation taken considering policy ϕ is played

Reinforcement Learning Setup



Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;

$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$ – reward function

★ ψ – Learning algorithm

★ \mathbb{E}_ϕ – Expectation taken considering policy ϕ is played

- Stationary deterministic policy: a function $\phi: \mathcal{S} \rightarrow \mathcal{A}$

Average reward criterion:

- Let

$$J_{\mathcal{M}}(\phi) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\phi \left[\sum_{t=1}^T r_t \right]$$

- Optimal average reward

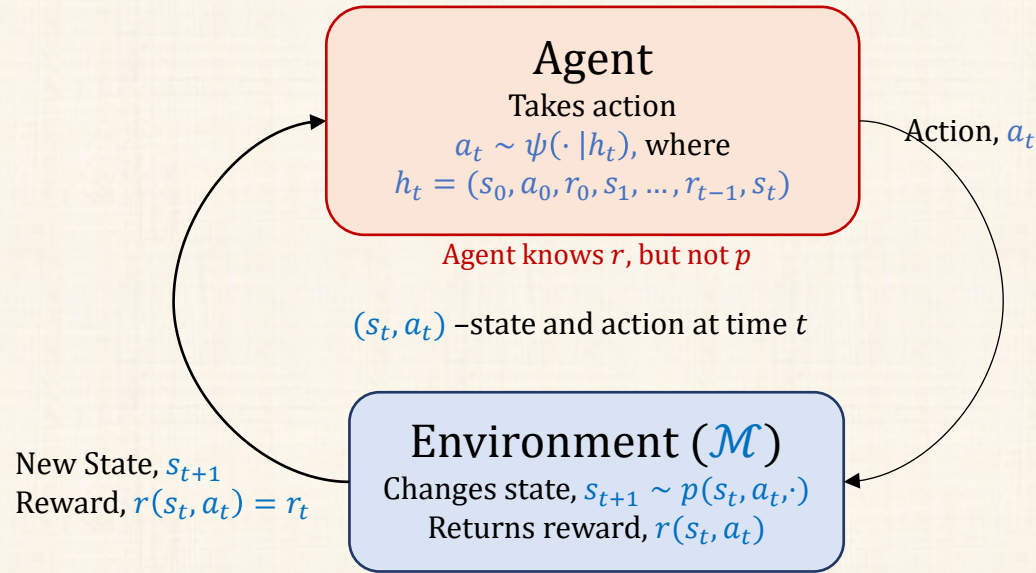
$$J_{\mathcal{M}}^* = \max_{\phi} J_{\mathcal{M}}(\phi)$$

Regret:

- Cumulative regret of algorithm, ψ is defined as

$$R(T; \psi) := TJ_{\mathcal{M}}^* - \sum_{t=1}^T r_t$$

Reinforcement Learning Setup



Model of the Environment: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$

- \mathcal{S} – state space, \mathcal{A} – action space
- p – transition kernel;

$$p(s, a, B) = \mathbb{P}(s_{t+1} \in B \mid s_t = s, a_t = a)$$
- $r: \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ – reward function

- Stationary deterministic policy: a function $\phi: \mathcal{S} \rightarrow \mathcal{A}$

Average reward criterion:

- Let

$$J_{\mathcal{M}}(\phi) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\phi} \left[\sum_{t=1}^T r_t \right]$$

- Optimal average reward

$$J_{\mathcal{M}}^* = \max_{\phi} J_{\mathcal{M}}(\phi)$$

Regret:

- Cumulative regret of algorithm, ψ is defined as

$$R(T; \psi) := TJ_{\mathcal{M}}^* - \sum_{t=1}^T r_t$$

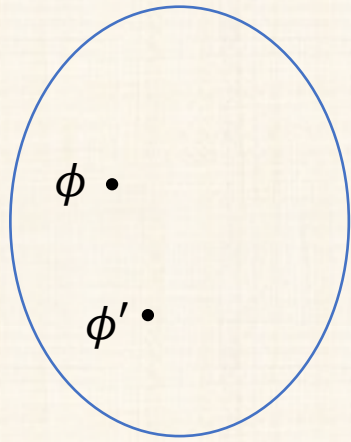
Goal: To propose ψ with low regret upper bound

★ ψ – Learning algorithm

★ \mathbb{E}_{ϕ} – Expectation taken considering policy ϕ is played

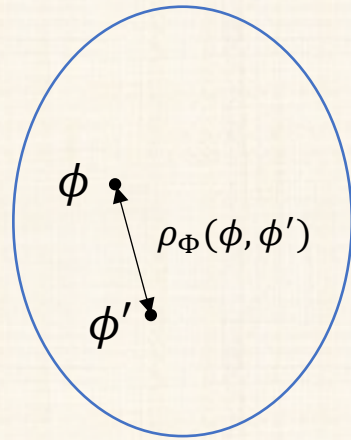
Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$



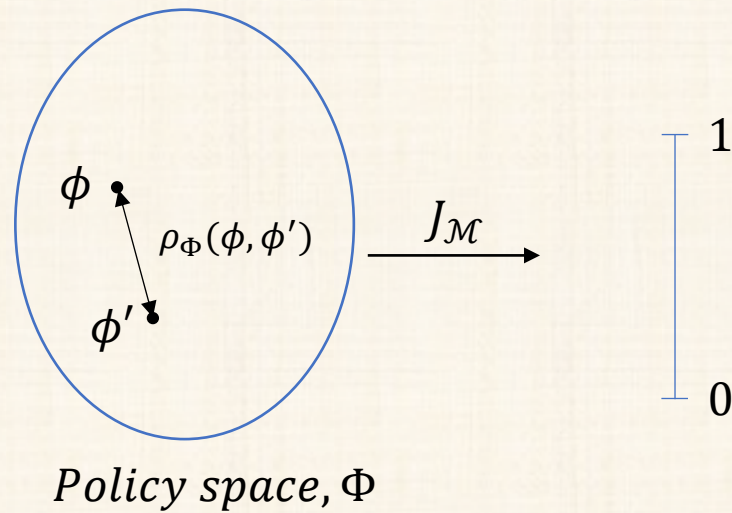
Policy space, Φ

Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

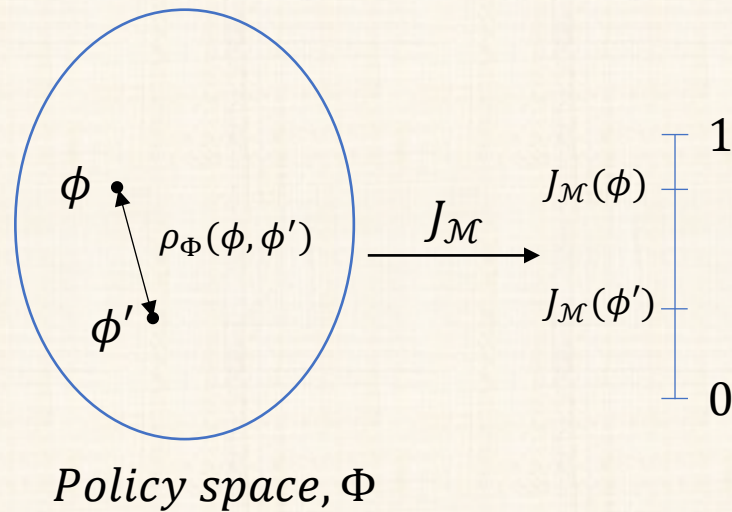


Policy space, Φ

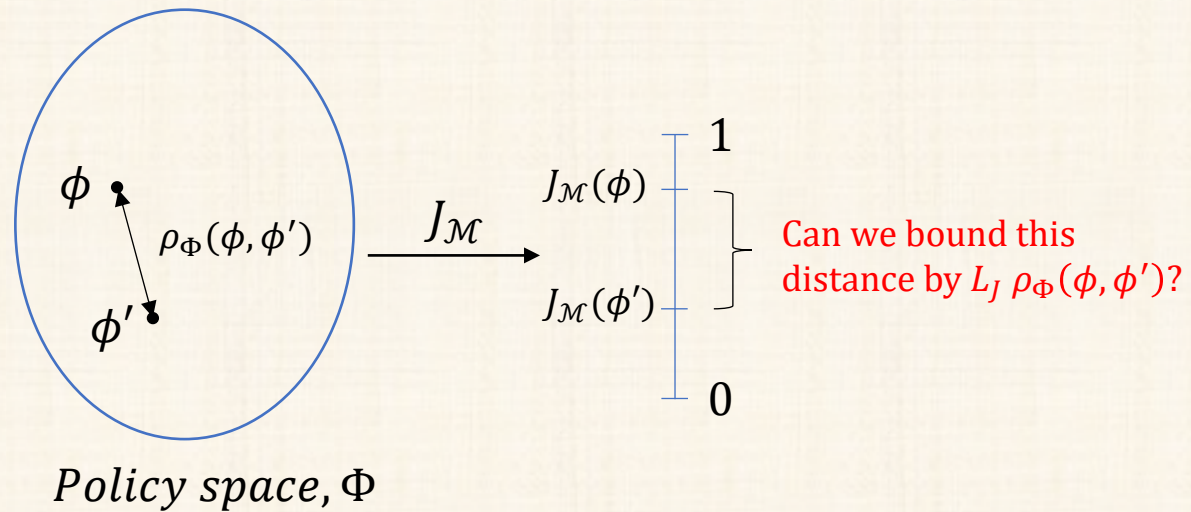
Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$



Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

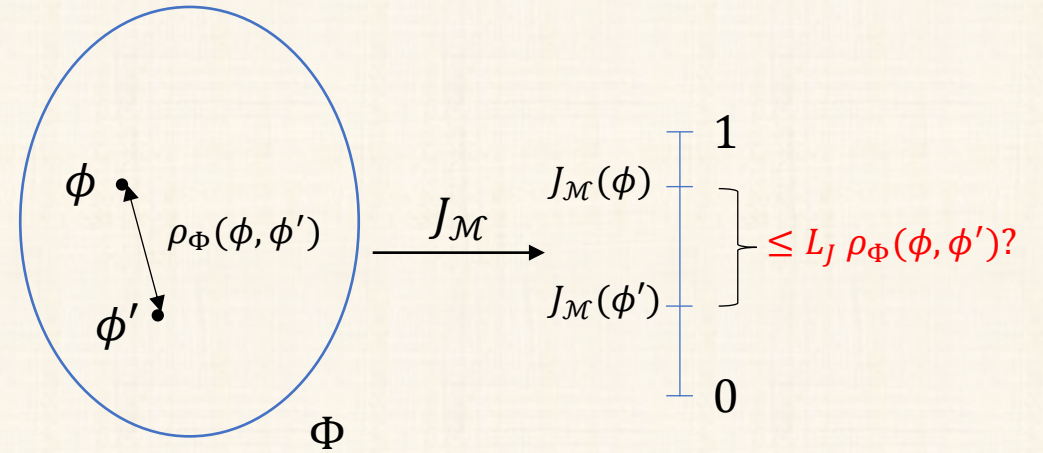


Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$



Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions



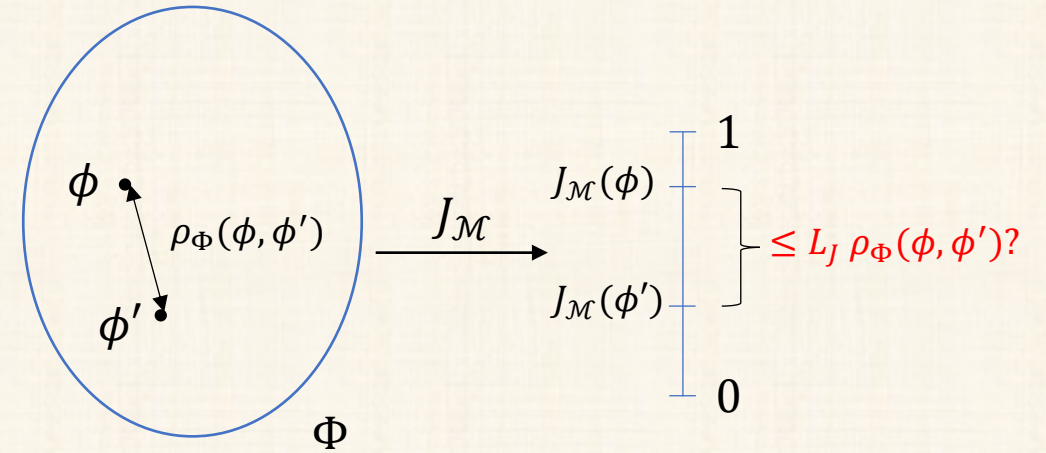
Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions

1. Lipschitz continuity: For every $(s, a), (s', a')$
 $|r(s, a) - r(s', a')| \leq L_r \rho((s, a), (s', a'))$,
 $\|p(s, a, \cdot) - p(s', a', \cdot)\|_{TV} \leq L_p \rho((s, a), (s', a'))$.
2. Uniform ergodicity: $\exists \alpha \in [0, 1)$ and $C < \infty$ such that for every $(s, a), (s', a')$

$$\left\| \mu_{\phi, p, s_0}^{(t)} - \mu_{\phi, p}^{(\infty)} \right\|_{TV} \leq C \cdot \alpha^t, \forall t \in \mathbb{N}.$$

3. Upper bound on stationary measure: $\exists \bar{\kappa} > 0$ and a probability measure ν such that for every $\phi, \mu_{\phi, p}^{(\infty)} \leq \bar{\kappa} \cdot \nu$.



* $\mu_{\phi, p, s'}^{(t)}$ - distribution of s_t when ϕ is played and initial state is s'

* $\mu_{\phi, p}^{(\infty)}$ - stationary distribution of $\{s_t\}$ under application of ϕ

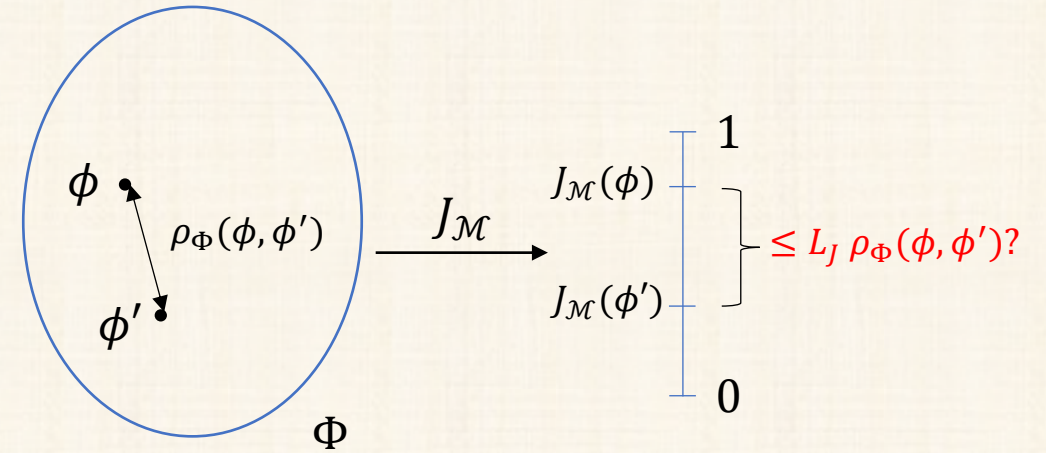
Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions

1. Lipschitz continuity: For every $(s, a), (s', a')$
 $|r(s, a) - r(s', a')| \leq L_r \rho((s, a), (s', a'))$,
 $\|p(s, a, \cdot) - p(s', a', \cdot)\|_{TV} \leq L_p \rho((s, a), (s', a'))$.
2. Uniform ergodicity: $\exists \alpha \in [0, 1)$ and $C < \infty$ such that for every $(s, a), (s', a')$

$$\left\| \mu_{\phi, p, s_0}^{(t)} - \mu_{\phi, p}^{(\infty)} \right\|_{TV} \leq C \cdot \alpha^t, \forall t \in \mathbb{N}.$$

3. Upper bound on stationary measure: $\exists \bar{\kappa} > 0$ and a probability measure ν such that for every ϕ , $\mu_{\phi, p}^{(\infty)} \leq \bar{\kappa} \cdot \nu$.



Metric, ρ_{Φ}

- * $\mu_{\phi, p, s'}^{(t)}$ - distribution of s_t when ϕ is played and initial state is s'
- * $\mu_{\phi, p}^{(\infty)}$ - stationary distribution of $\{s_t\}$ under application of ϕ

Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions

1. Lipschitz continuity: For every $(s, a), (s', a')$
 $|r(s, a) - r(s', a')| \leq L_r \rho((s, a), (s', a')),$
 $\|p(s, a, \cdot) - p(s', a', \cdot)\|_{TV} \leq L_p \rho((s, a), (s', a')).$
2. Uniform ergodicity: $\exists \alpha \in [0, 1)$ and $C < \infty$ such that for every $(s, a), (s', a')$

$$\left\| \mu_{\phi, p, s_0}^{(t)} - \mu_{\phi, p}^{(\infty)} \right\|_{TV} \leq C \cdot \alpha^t, \forall t \in \mathbb{N}.$$

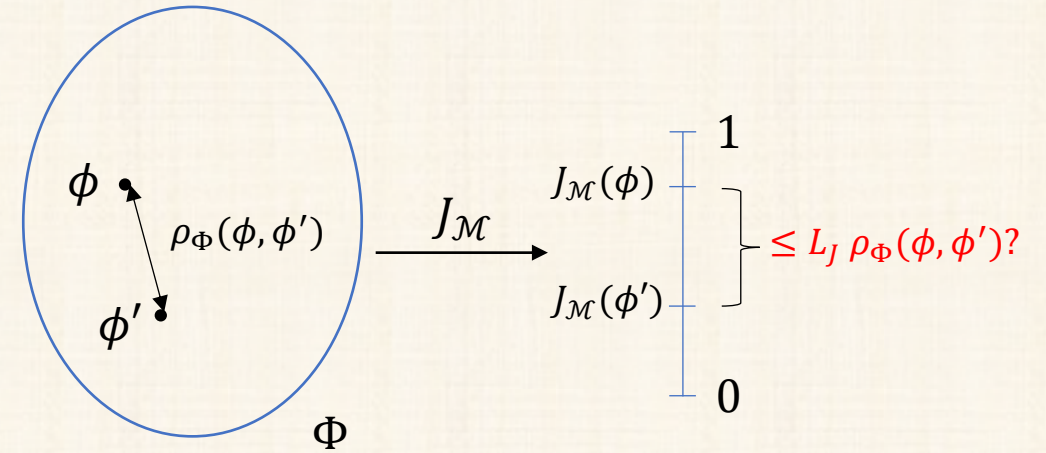
3. Upper bound on stationary measure: $\exists \bar{\kappa} > 0$ and a probability measure ν such that for every $\phi, \mu_{\phi, p}^{(\infty)} \leq \bar{\kappa} \cdot \nu$.

Metric, ρ_{Φ}

- $\rho_{\Phi}(\phi, \phi') := \int \rho_{\mathcal{A}}(\phi(s), \phi'(s)) d\nu(s)$, where ν satisfies 3 and $\rho_{\mathcal{A}}$ is a metric on \mathcal{A} .

★ $\mu_{\phi, p, s'}^{(t)}$ - distribution of s_t when ϕ is played and initial state is s'

★ $\mu_{\phi, p}^{(\infty)}$ - stationary distribution of $\{s_t\}$ under application of ϕ



Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions

1. Lipschitz continuity: For every $(s, a), (s', a')$
 $|r(s, a) - r(s', a')| \leq L_r \rho((s, a), (s', a')),$
 $\|p(s, a, \cdot) - p(s', a', \cdot)\|_{TV} \leq L_p \rho((s, a), (s', a')).$
2. Uniform ergodicity: $\exists \alpha \in [0, 1)$ and $C < \infty$ such that for every $(s, a), (s', a')$

$$\left\| \mu_{\phi, p, s_0}^{(t)} - \mu_{\phi, p}^{(\infty)} \right\|_{TV} \leq C \cdot \alpha^t, \forall t \in \mathbb{N}.$$

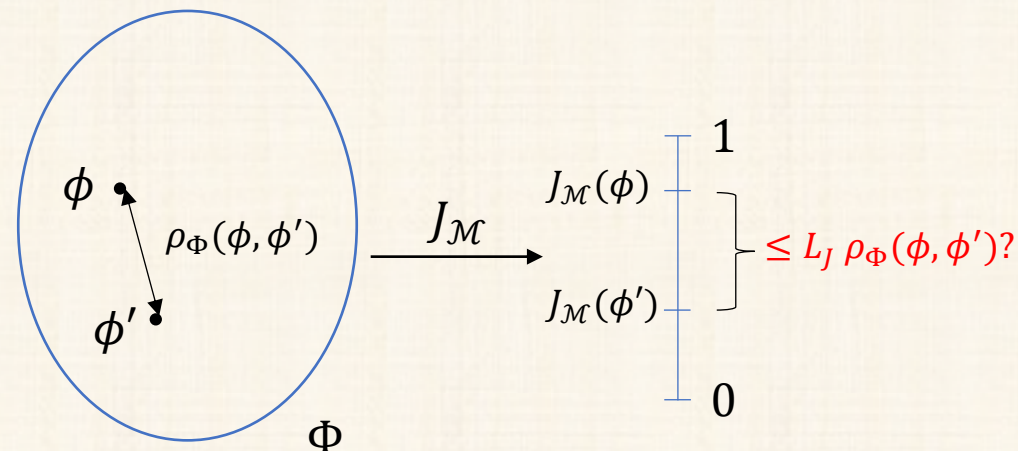
3. Upper bound on stationary measure: $\exists \bar{\kappa} > 0$ and a probability measure ν such that for every $\phi, \mu_{\phi, p}^{(\infty)} \leq \bar{\kappa} \cdot \nu$.

Metric, ρ_{Φ}

- $\rho_{\Phi}(\phi, \phi') := \int \rho_{\mathcal{A}}(\phi(s), \phi'(s)) d\nu(s)$, where ν satisfies 3 and $\rho_{\mathcal{A}}$ is a metric on \mathcal{A} .

★ $\mu_{\phi, p, s'}^{(t)}$ - distribution of s_t when ϕ is played and initial state is s'

★ $\mu_{\phi, p}^{(\infty)}$ - stationary distribution of $\{s_t\}$ under application of ϕ



Theorem: Under Assumption 1, 2 and 3, for any ϕ, ϕ'

$$\left\| \mu_{\phi, p}^{(\infty)} - \mu_{\phi', p}^{(\infty)} \right\|_{TV} \leq (\lceil \log_{\alpha^{-1}}(C) \rceil + 1) \frac{\bar{\kappa} L_p}{1 - \alpha} \rho_{\Phi}(\phi, \phi')$$

Lipschitz Continuity of $J_{\mathcal{M}}(\cdot)$

Assumptions

1. Lipschitz continuity: For every $(s, a), (s', a')$
 $|r(s, a) - r(s', a')| \leq L_r \rho((s, a), (s', a')),$
 $\|p(s, a, \cdot) - p(s', a', \cdot)\|_{TV} \leq L_p \rho((s, a), (s', a')).$
2. Uniform ergodicity: $\exists \alpha \in [0, 1)$ and $C < \infty$ such that for every $(s, a), (s', a')$

$$\left\| \mu_{\phi, p, s_0}^{(t)} - \mu_{\phi, p}^{(\infty)} \right\|_{TV} \leq C \cdot \alpha^t, \forall t \in \mathbb{N}.$$

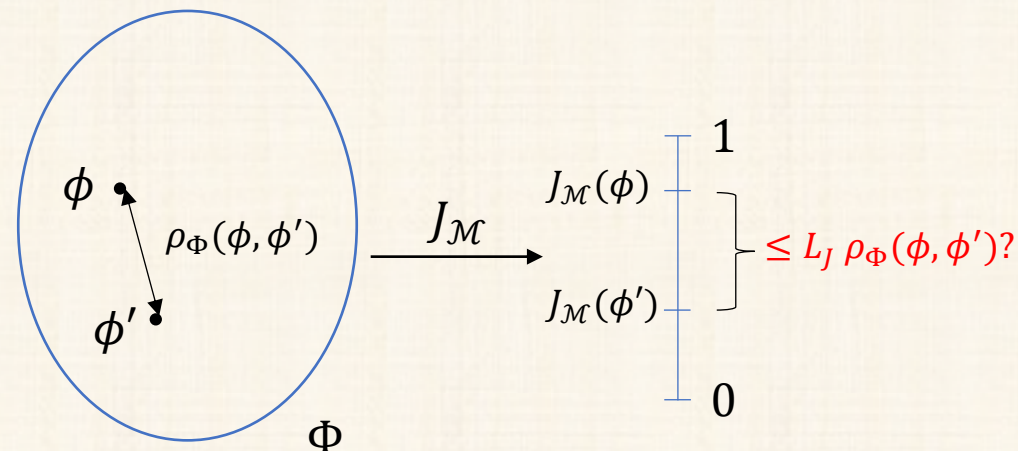
3. Upper bound on stationary measure: $\exists \bar{\kappa} > 0$ and a probability measure ν such that for every $\phi, \mu_{\phi, p}^{(\infty)} \leq \bar{\kappa} \cdot \nu$.

Metric, ρ_{Φ}

- $\rho_{\Phi}(\phi, \phi') := \int \rho_{\mathcal{A}}(\phi(s), \phi'(s)) d\nu(s)$, where ν satisfies 3 and $\rho_{\mathcal{A}}$ is a metric on \mathcal{A} .

★ $\mu_{\phi, p, s'}^{(t)}$ - distribution of s_t when ϕ is played and initial state is s'

★ $\mu_{\phi, p}^{(\infty)}$ - stationary distribution of $\{s_t\}$ under application of ϕ



Theorem: Under Assumption 1, 2 and 3, for any ϕ, ϕ'

$$\left\| \mu_{\phi, p}^{(\infty)} - \mu_{\phi', p}^{(\infty)} \right\|_{TV} \leq (\lceil \log_{\alpha^{-1}}(C) \rceil + 1) \frac{\bar{\kappa} L_p}{1 - \alpha} \rho_{\Phi}(\phi, \phi')$$

Corollary: Under Assumption 1, 2 and 3, for any ϕ, ϕ'

$$|J_{\mathcal{M}}(\phi) - J_{\mathcal{M}}(\phi')| \leq L_J \rho_{\Phi}(\phi, \phi')$$

where $L_J := \bar{\kappa} \left(L_r + \frac{(\lceil \log_{\alpha^{-1}}(C) \rceil + 1) L_p}{(1 - \alpha)} \right)$.

Zooming in Policy Space

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{\mathcal{O}}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$

Compact policy set, Φ

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{\mathcal{O}}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

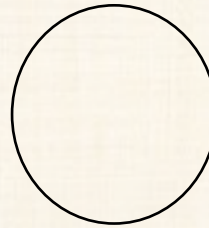
$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{\mathcal{O}}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ

Compact policy set, Φ

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{\mathcal{O}}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

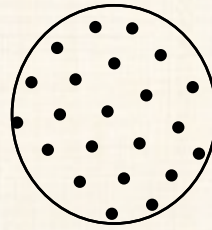
$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{\mathcal{O}}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ

Compact policy set, Φ

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$

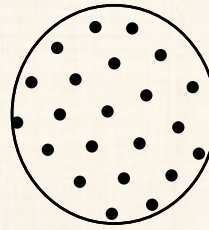
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

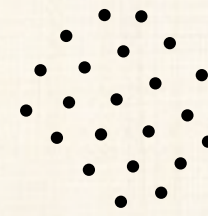
$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ



ϵ -net of Φ , Φ_{ϵ}

Compact policy set, Φ

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$

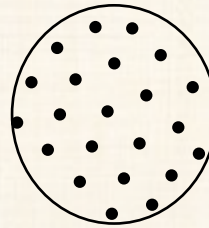
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

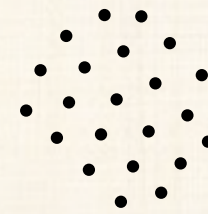
$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ



ϵ -net of Φ , Φ_{ϵ}

Compact policy set, Φ

If

1. $J_{\mathcal{M}}: \Phi \rightarrow [0,1]$ is L_J -Lipschitz,
2. The algorithm for finite policy set is run with Φ_{ϵ}
3. $\epsilon = T^{-\frac{\beta}{d^{\Phi}\beta+1}}$, where d^{Φ} is the dimension Φ ,

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$

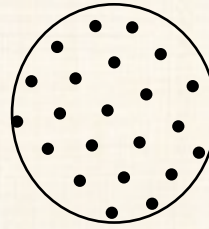
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

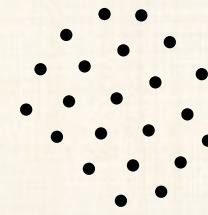
$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ



ϵ -net of Φ , Φ_{ϵ}

Compact policy set, Φ

If

1. $J_{\mathcal{M}}: \Phi \rightarrow [0,1]$ is L_J -Lipschitz,
2. The algorithm for finite policy set is run with Φ_{ϵ}
3. $\epsilon = T^{-\frac{\beta}{d^{\Phi}\beta+1}}$, where d^{Φ} is the dimension Φ ,

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}\left(T^{\frac{(d^{\Phi}-1)\beta+1}{d^{\Phi}\beta+1}}\right)$$

Zooming in Policy Space

Define Regret w.r.t. Φ as,

$$R_{\Phi}(T; \psi) := T \max_{\phi \in \Phi} J_{\mathcal{M}}(\phi) - \sum_{t=1}^T r_t$$

Finite policy set, Φ

Suppose that we have

$$|J_{\mathcal{M}}(\phi) - \bar{r}_t(\phi)| \leq c_t(\phi) = \tilde{O}\left(\frac{1}{N_t(\phi)^{\beta}}\right)$$

where

$$N_t(\phi) = \sum_{s < t} \mathbb{I}(\phi_s = \phi), \text{ and}$$

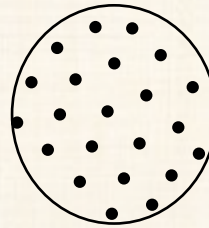
$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_s r_s \mathbb{I}(\phi_s = \phi)$$

Choose $\phi_t \in \operatorname{argmax}_{\phi \in \Phi} \operatorname{Index}_t(\phi)$, where

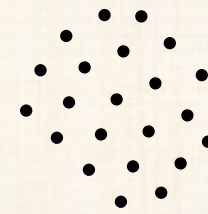
$$\operatorname{Index}_t(\phi) := \bar{r}_t(\phi) + \frac{\text{const}}{N_t(\phi)^{\beta}}$$

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}(|\Phi|^{\beta} T^{1-\beta})$$



Policy space, Φ



ϵ -net of Φ , Φ_{ϵ}

Compact policy set, Φ

If

1. $J_{\mathcal{M}}: \Phi \rightarrow [0,1]$ is L_J -Lipschitz,
2. The algorithm for finite policy set is run with Φ_{ϵ}
3. $\epsilon = T^{-\frac{\beta}{d^{\Phi}\beta+1}}$, where d^{Φ} is the dimension Φ ,

Then,

$$R_{\Phi}(T; \psi) \leq \tilde{O}\left(T^{\frac{(d^{\Phi}-1)\beta+1}{d^{\Phi}\beta+1}}\right)$$

- For example, $\beta = 0.5 \Rightarrow R_{\Phi}(T; \psi) \leq \tilde{O}\left(T^{\frac{d^{\Phi}+1}{d^{\Phi}+2}}\right)$

Zooming in Policy Space

Policy Zooming for RL

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$,
 $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Zooming in Policy Space

Policy Zooming for RL

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$,
 $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Model-free

Model-based

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$,
 $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_{\tau} r_{\tau} \mathbb{I}(\phi_{\tau} = \phi),$$

$$c_t(\phi) = \text{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_{\tau} \mathbb{I}(\phi_{\tau} = \phi)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Zooming dimension:

$$d_Z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_Z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

$$c_t(\phi) = \text{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

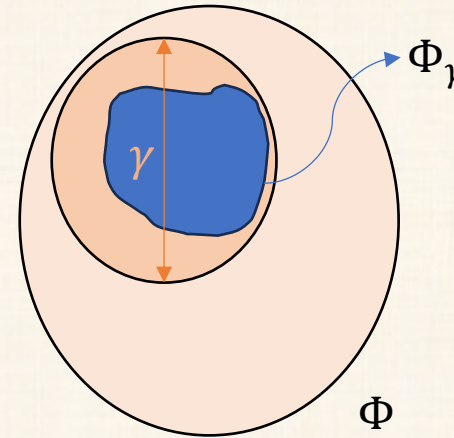
$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Zooming dimension:

$$d_z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

$$c_t(\phi) = \text{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

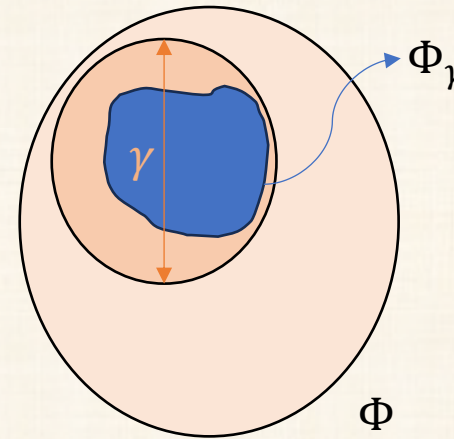
$$c_t(\phi) = \text{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming dimension:

$$d_Z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_Z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



In this case, $d_Z^\Phi = 1$, but $d^\Phi = 2$.

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

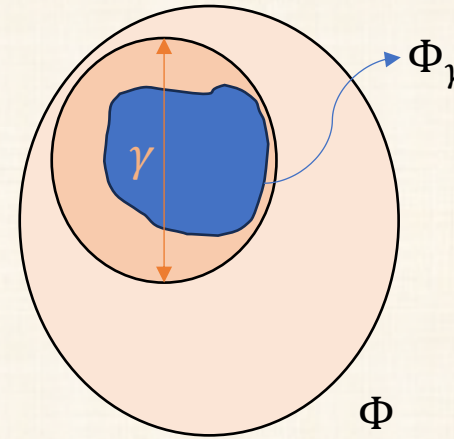
$$c_t(\phi) = \operatorname{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming dimension:

$$d_z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



In this case, $d_z^\Phi = 1$, but $d^\Phi = 2$.

Theorem: If \mathcal{M} satisfies Assumption 1, 2 and 3, then with high probability

$$R_\Phi(T; \text{PZRLMF}) \leq \tilde{O}\left(T^{\frac{d_z^\Phi + 1}{d_z^\Phi + 2}}\right)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

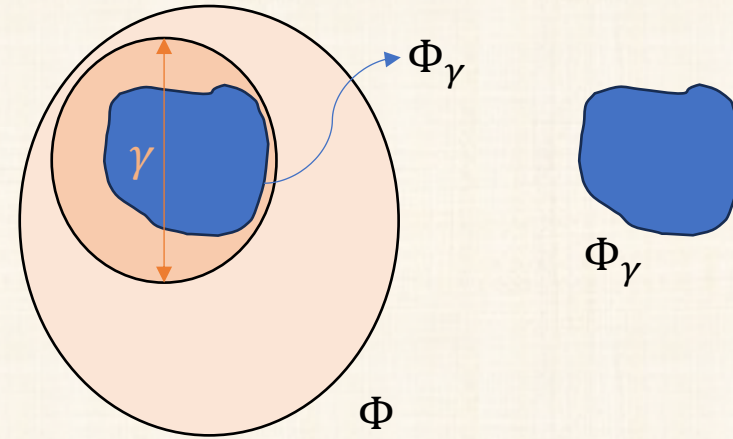
$$c_t(\phi) = \operatorname{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming dimension:

$$d_z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



In this case, $d_z^\Phi = 1$, but $d^\Phi = 2$.

Theorem: If \mathcal{M} satisfies Assumption 1, 2 and 3, then with high probability

$$R_\Phi(T; \text{PZRLMF}) \leq \tilde{O} \left(T^{\frac{d_z^\Phi + 1}{d_z^\Phi + 2}} \right)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

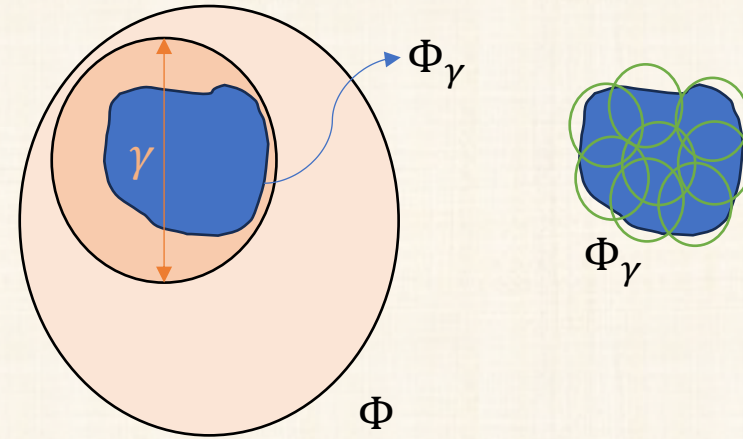
$$c_t(\phi) = \operatorname{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming dimension:

$$d_z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



In this case, $d_z^\Phi = 1$, but $d^\Phi = 2$.

Theorem: If \mathcal{M} satisfies Assumption 1, 2 and 3, then with high probability

$$R_\Phi(T; \text{PZRLMF}) \leq \tilde{O} \left(T^{\frac{d_z^\Phi + 1}{d_z^\Phi + 2}} \right)$$

Zooming in Policy Space

Policy Zooming for RL (Model-free)

Inputs: Horizon T , A policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

where

$$\operatorname{Index}_t(\phi) = \bar{r}_t(\phi) + c_t(\phi)$$

$$\bar{r}_t(\phi) = \frac{1}{N_t(\phi)} \sum_t r_t \mathbb{I}(\phi_t = \phi),$$

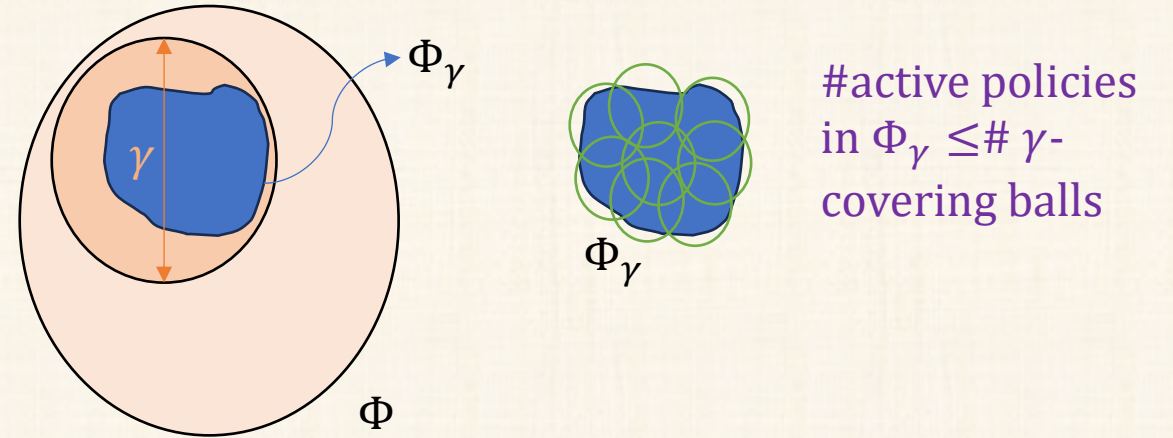
$$c_t(\phi) = \operatorname{const} \cdot N_t(\phi)^{-\frac{1}{2}}, \text{ and}$$

$$N_t(\phi) = \sum_t \mathbb{I}(\phi_t = \phi)$$

Zooming dimension:

$$d_z^\Phi = \inf\{d > 0: N_\gamma(\Phi_\gamma) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where Φ_γ is the set of $[\gamma, 2\gamma)$ -suboptimal policies.



In this case, $d_z^\Phi = 1$, but $d^\Phi = 2$.

Theorem: If \mathcal{M} satisfies Assumption 1, 2 and 3, then with high probability

$$R_\Phi(T; \text{PZRLMF}) \leq \tilde{O} \left(T^{\frac{d_z^\Phi + 1}{d_z^\Phi + 2}} \right)$$

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$,
 $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

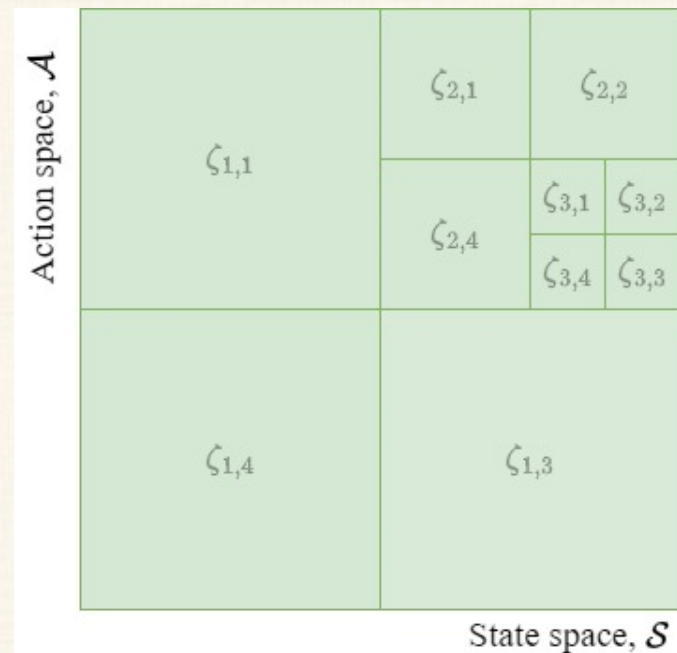
$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

1. Adaptively discretized partition of $\mathcal{S} \times \mathcal{A}$:

ζ is active at time $t \Leftrightarrow N_t(\zeta) \geq \operatorname{diam}(\zeta)^{-(d_s+2)}$ but $< 2^{d_s+2} \operatorname{diam}(\zeta)^{-(d_s+2)}$



Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

1. Adaptively discretized partition of $\mathcal{S} \times \mathcal{A}$:

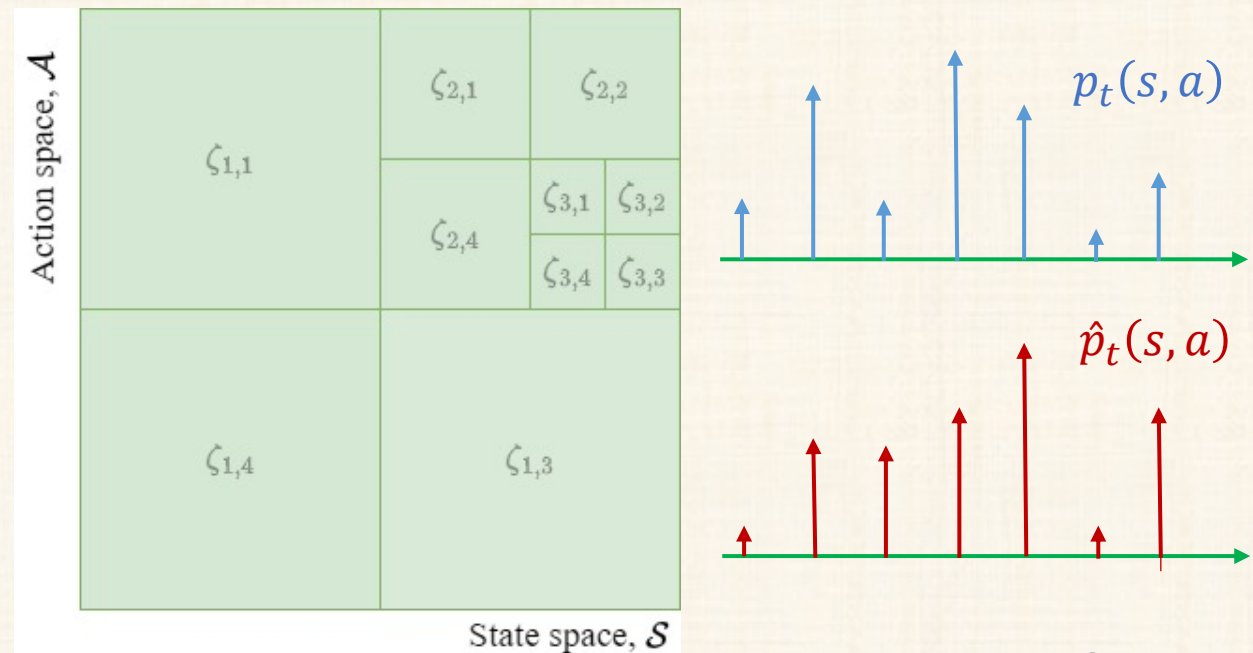
ζ is active at time $t \Leftrightarrow N_t(\zeta) \geq \operatorname{diam}(\zeta)^{-(d_s+2)}$ but $< 2^{d_s+2} \operatorname{diam}(\zeta)^{-(d_s+2)}$

2. Confidence set for p :

• Z_t : Representative point of cells

$\mathcal{C}_t = \{\theta: \|\theta(s, a) - \hat{p}_t(s, a)\|_{TV} \leq \operatorname{diam}(\zeta_{s,a}) \forall (s, a) \in Z_t\}$

Lemma: $p_t \in \mathcal{C}_t \forall t \in [T]$ with high probability.



Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

1. Adaptively discretized partition of $\mathcal{S} \times \mathcal{A}$:

ζ is active at time $t \Leftrightarrow N_t(\zeta) \geq \operatorname{diam}(\zeta)^{-(d_S+2)}$ but $< 2^{d_S+2} \operatorname{diam}(\zeta)^{-(d_S+2)}$

2. Confidence set for p :

• Z_t : Representative point of cells

$\mathcal{C}_t = \{\theta: \|\theta(s, a) - \hat{p}_t(s, a)\|_{TV} \leq \operatorname{diam}(\zeta_{s,a}) \forall (s, a) \in Z_t\}$

Lemma: $p_t \in \mathcal{C}_t \forall t \in [T]$ with high probability.

3. Index of policy ϕ :

$\bar{V}_0^\phi(s) = 0$

$\bar{V}_{i+1}^\phi(s) = r^+(\zeta_{s,\phi(s)}) + \max_{\theta \in \mathcal{C}_t} \sum \theta(q(\zeta_{s,\phi(s)}), s') \bar{V}_i^\phi(s')$

where $r^+ = r + \operatorname{const} \cdot \operatorname{diam}$,

$q(\cdot)$ denotes quantized point

$\operatorname{Index}_t(\phi) := \lim_{i \rightarrow \infty} \frac{1}{i} \bar{V}_i^\phi(s)$

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Assumptions (Contd.)

4. Lower bound on stationary measure: $\exists \underline{\kappa} > 0$ such that for every ϕ , $\mu_{\phi,p}^{(\infty)} \geq \underline{\kappa} \cdot \lambda$, where λ is the Lebesgue measure on \mathcal{S} .
5. Partial derivatives of transition densities are bounded.

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Assumptions (Contd.)

4. Lower bound on stationary measure: $\exists \underline{\kappa} > 0$ such that for every ϕ , $\mu_{\phi,p}^{(\infty)} \geq \underline{\kappa} \cdot \lambda$, where λ is the Lebesgue measure on \mathcal{S} .
5. Partial derivatives of transition densities are bounded.

Lemma: Under Assumption 1 and 5, with high probability $J_{\mathcal{M}}(\phi) \leq \operatorname{Index}_t(\phi), \forall t = 1, 2, \dots, T$.

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{\}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \bigcup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Assumptions (Contd.)

4. Lower bound on stationary measure: $\exists \underline{\kappa} > 0$ such that for every ϕ , $\mu_{\phi,p}^{(\infty)} \geq \underline{\kappa} \cdot \lambda$, where λ is the Lebesgue measure on \mathcal{S} .
5. Partial derivatives of transition densities are bounded.

Lemma: Under Assumption 1 and 5, with high probability $J_{\mathcal{M}}(\phi) \leq \operatorname{Index}_t(\phi), \forall t = 1, 2, \dots, T$.

Lemma: Under Assumption 1, 2, 3, 4 and 5, with high probability

$$\begin{aligned} & \operatorname{Index}_t(\phi) \\ & \leq J_{\mathcal{M}}(\phi) + \text{const} \cdot \int \operatorname{diam}(\zeta_{s,\phi(s)}) d\mu_{\phi,p}^{(\infty)}(s), \\ & \quad \forall t = 1, 2, \dots, T. \end{aligned}$$

Zooming in Policy Space

Policy Zooming for RL (Model-based)

Inputs: Horizon T , A Lipschitz policy class Φ

Initialize: Set of active policies $\Phi_0 \leftarrow \{ \}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

$\Phi_t = \Phi_{t-1}$

while $\exists \phi' \in \Phi$ such that $\phi' \in \cup_{\phi \in \Phi_t} B(\phi; c_t(\phi))$:

$\Phi_t \leftarrow \Phi_t \cup \{\phi'\}$

$\phi_k \in \operatorname{argmax}_{\phi \in \Phi_t} \operatorname{Index}_t(\phi)$

$H = \max\{1, N_t(\phi_k)\}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Assumptions (Contd.)

4. Lower bound on stationary measure: $\exists \underline{\kappa} > 0$ such that for every $\phi, \mu_{\phi,p}^{(\infty)} \geq \underline{\kappa} \cdot \lambda$, where λ is the Lebesgue measure on \mathcal{S} .
5. Partial derivatives of transition densities are bounded.

Lemma: Under Assumption 1 and 5, with high probability $J_{\mathcal{M}}(\phi) \leq \operatorname{Index}_t(\phi), \forall t = 1, 2, \dots, T$.

Lemma: Under Assumption 1, 2, 3, 4 and 5, with high probability

$$\begin{aligned} & \operatorname{Index}_t(\phi) \\ & \leq J_{\mathcal{M}}(\phi) + \text{const} \cdot \int \operatorname{diam}(\zeta_{s,\phi(s)}) d\mu_{\phi,p}^{(\infty)}(s), \\ & \quad \forall t = 1, 2, \dots, T. \end{aligned}$$

Theorem: If \mathcal{M} satisfies Assumption 1, 2, 3, 4 and 5, then with high probability

$$R_{\Phi}(T; \text{PZRLMB}) \leq \tilde{O} \left(\frac{2d_{\mathcal{S}} + d_{\mathcal{Z}}^{\Phi} + 2}{T^{2d_{\mathcal{S}} + d_{\mathcal{Z}}^{\Phi} + 3}} \right)$$

Zooming in Policy Space

Two Special Cases:

Parameterized policy space:

- $\phi(\cdot; w)$ – policy parameterized by w
- $w \in W \subset \mathbb{R}^{d_w}$
- $\|\phi(\cdot; w_1) - \phi(\cdot; w_2)\|_v \leq \|w_1 - w_2\|$

Then

$$d_z^\Phi \leq d_w$$

Zooming in Policy Space

Two Special Cases:

Parameterized policy space:

- $\phi(\cdot; w)$ – policy parameterized by w
- $w \in W \subset \mathbb{R}^{d_w}$
- $\|\phi(\cdot; w_1) - \phi(\cdot; w_2)\|_v \leq \|w_1 - w_2\|$

Then

$$d_Z^\Phi \leq d_w$$

Curvature condition:

If there is a unique maximum of $J_{\mathcal{M}}(\phi(\cdot; \cdot)): W \rightarrow [0,1]$ at w^* , and

$$J_{\mathcal{M}}(\phi(\cdot; w^*)) - J_{\mathcal{M}}(\phi(\cdot; w)) \geq K_w \|w - w^*\|, \forall w \in W.$$

Then

$$d_Z^\Phi = 0$$

Zooming in Policy Space

Two Special Cases:

Parameterized policy space:

- $\phi(\cdot; w)$ – policy parameterized by w
- $w \in W \subset \mathbb{R}^{d_w}$
- $\|\phi(\cdot; w_1) - \phi(\cdot; w_2)\|_v \leq \|w_1 - w_2\|$

Then

$$d_z^\Phi \leq d_w$$

Curvature condition:

If there is a unique maximum of $J_{\mathcal{M}}(\phi(\cdot; \cdot)): W \rightarrow [0,1]$ at w^* , and

$$J_{\mathcal{M}}(\phi(\cdot; w^*)) - J_{\mathcal{M}}(\phi(\cdot; w)) \geq K_w \|w - w^*\|, \forall w \in W.$$

Then

$$d_z^\Phi = 0$$

Drawbacks:

- Dimension of Φ could be huge
- Knowledge of low-complexity Φ may not be available
- Computationally heavy

Zooming in State-Action Space

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0, k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0$, $k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

$\star \exists h_{\mathcal{M}}: \mathcal{S} \rightarrow \mathbb{R}$ s.t.

$$J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s')\}$$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0$, $k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

★ $\exists h_{\mathcal{M}}: \mathcal{S} \rightarrow \mathbb{R}$ s.t.

$$J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s')\}$$

$$\star \text{gap}(s, a) := J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) - (r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s'))$$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0$, $k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

★ $\exists h_{\mathcal{M}}: \mathcal{S} \rightarrow \mathbb{R}$ s.t.

$$J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s')\}$$

$$\star \text{gap}(s, a) := J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) - (r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s'))$$

Zooming dimension:

$$d_z = \inf\{d > 0: N_{\gamma}(Z_{\gamma}) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where $Z_{\gamma} = \{(s, a) \mid \text{gap}(s, a) \leq \gamma\}$.

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0$, $k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_{\mathcal{S}}+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

★ $\exists h_{\mathcal{M}}: \mathcal{S} \rightarrow \mathbb{R}$ s.t.

$$J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s')\}$$

$$\star \text{gap}(s, a) := J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) - (r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s'))$$

Zooming dimension:

$$d_z = \inf\{d > 0: N_{\gamma}(Z_{\gamma}) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

where $Z_{\gamma} = \{(s, a) \mid \text{gap}(s, a) \leq \gamma\}$.

$$\star d_z \leq d_{\mathcal{S}} + d_{\mathcal{A}}$$

Zooming in State-Action Space

Zooming for RL (ZoRL)

Inputs: Horizon T

Initialize: Set of $\mathcal{P}_0 = \mathcal{S} \times \mathcal{A}$, $k = 0$, $h = 0$, $H_0 = 0$

for $t = 1$ to T :

if $h = H_k$:

$h = 0$, $k \leftarrow k + 1$

 Update partition, \mathcal{P}_t

 Construct extended MDP, \mathcal{M}_t^+

$\phi_k \leftarrow \text{EVI}(\mathcal{M}_t^+, 1/\sqrt{T})$

$d_k \leftarrow \text{EPE}(\mathcal{M}_t^{d,+}, \phi_k, 1/\sqrt{T})$

$H_k = \text{const} \cdot d_k^{-2(d_S+1)}$

$h \leftarrow h + 1$

 Play $a_t = \phi_k(s_t)$

- $\mathcal{M}_t^+ = (S_t, A_t, \mathcal{C}_t, r_t^+)$ where
 $r_t^+(s, a) = r(s, a) + L_r \text{diam}(\zeta_{s,a}^t)$
- $\mathcal{M}_t^{d,+} = (S_t, A_t, \mathcal{C}_t, d_t)$ where
 $d_t(s, a) = \text{diam}(\zeta_{s,a}^t)$

Consider the average reward optimality equation (AROE) of \mathcal{M}

$$J + h(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h(s')\}$$

★ $\exists h_{\mathcal{M}}: \mathcal{S} \rightarrow \mathbb{R}$ s.t.

$$J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s')\}$$

$$\star \text{gap}(s, a) := J_{\mathcal{M}}^* + h_{\mathcal{M}}(s) - (r(s, a) + \int p(s, a, ds') h_{\mathcal{M}}(s'))$$

Zooming dimension:

$$d_z = \inf\{d > 0: N_{\gamma}(Z_{\gamma}) \leq c_z \gamma^{-d} \forall \gamma > 0\}$$

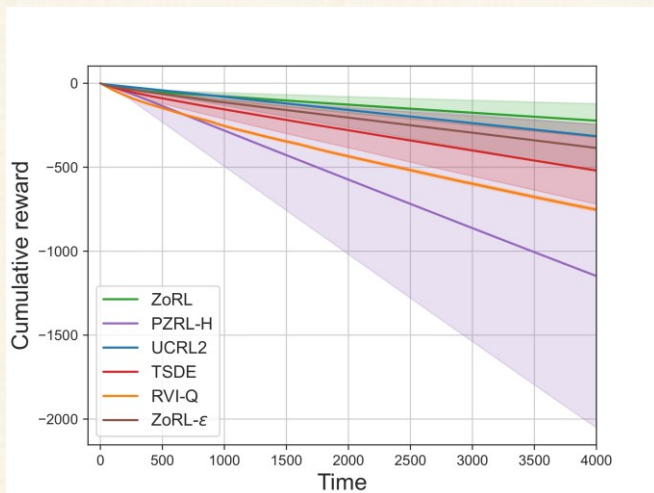
where $Z_{\gamma} = \{(s, a) \mid \text{gap}(s, a) \leq \gamma\}$.

$$\star d_z \leq d_S + d_{\mathcal{A}}$$

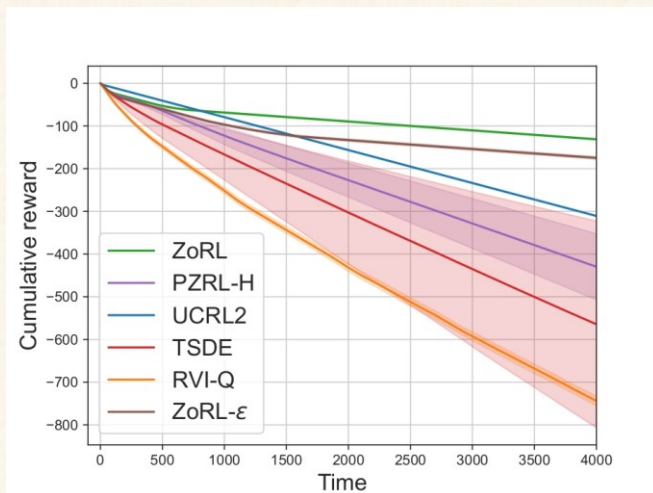
Theorem: If \mathcal{M} satisfies Assumption 1, 2, 3, 4 and 5, then with high probability

$$R(T; \text{ZoRL}) \leq \tilde{O}\left(T^{\frac{2d_S + d_z + 2}{2d_S + d_z + 3}}\right)$$

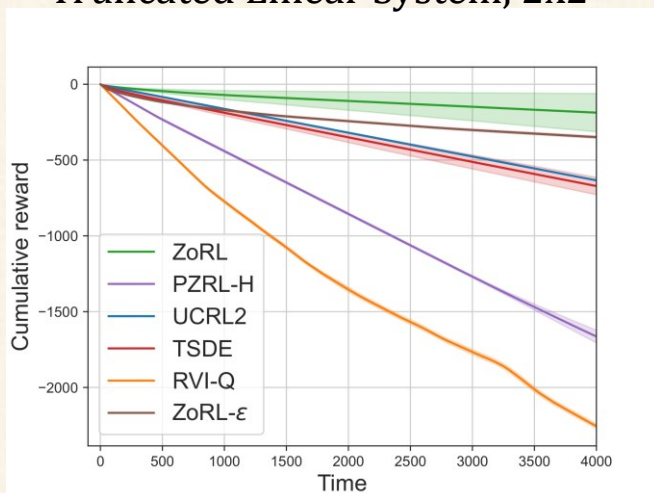
Zooming in State-Action Space



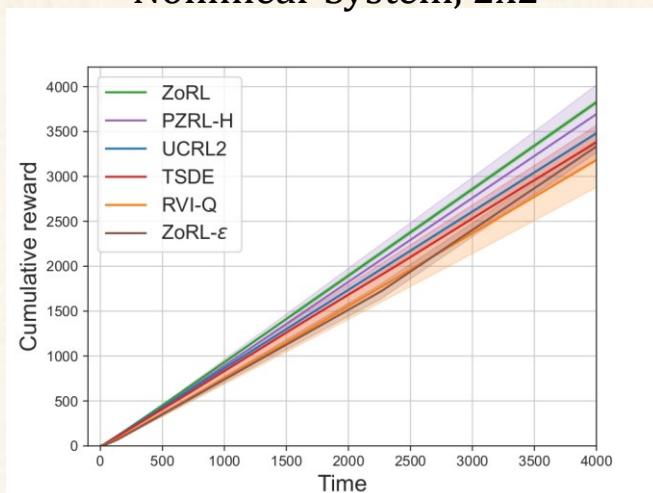
Truncated Linear System, 2x2



Nonlinear System, 2x2

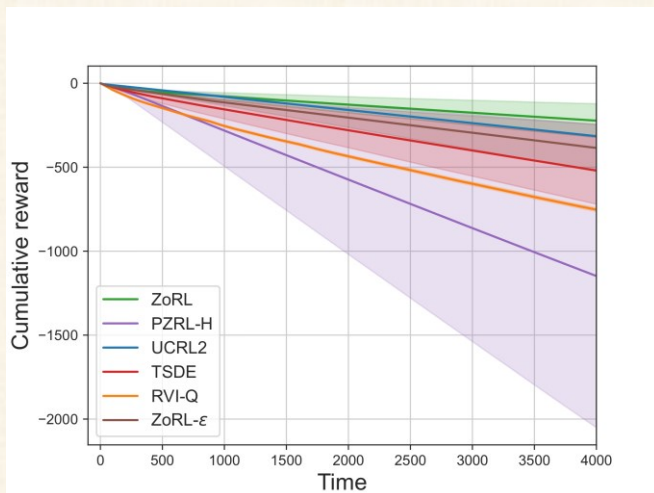


Truncated Linear System, 2x4

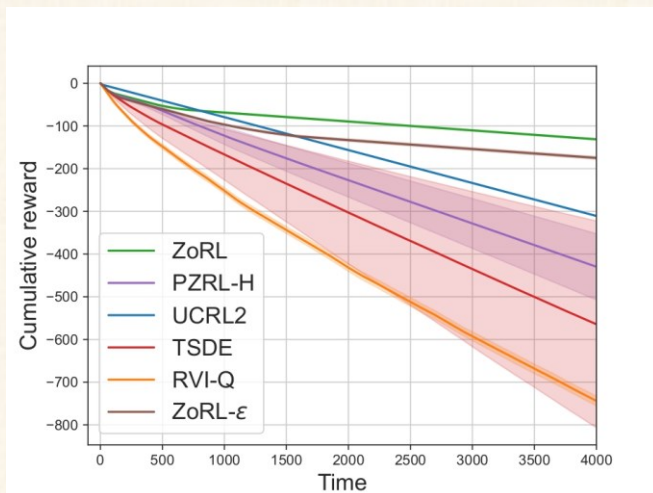


Continuous RiverSwim

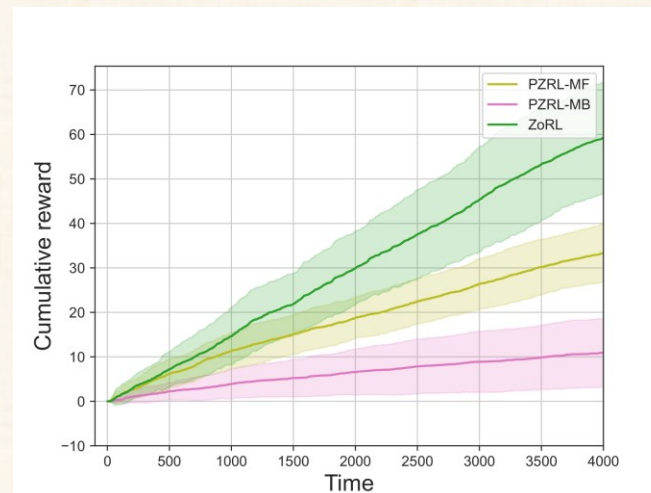
Zooming in State-Action Space



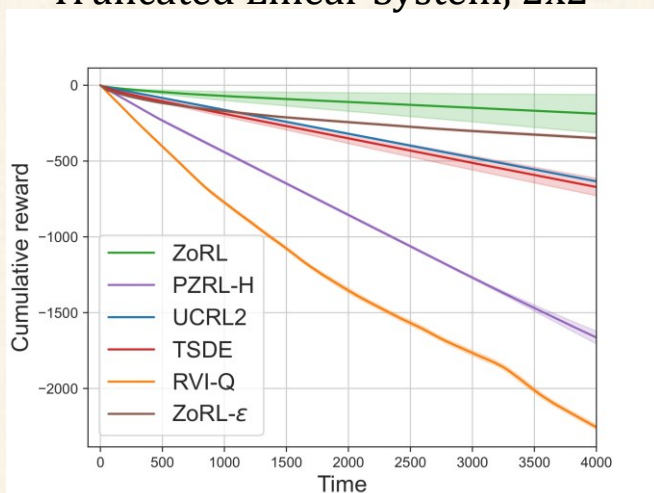
Truncated Linear System, 2x2



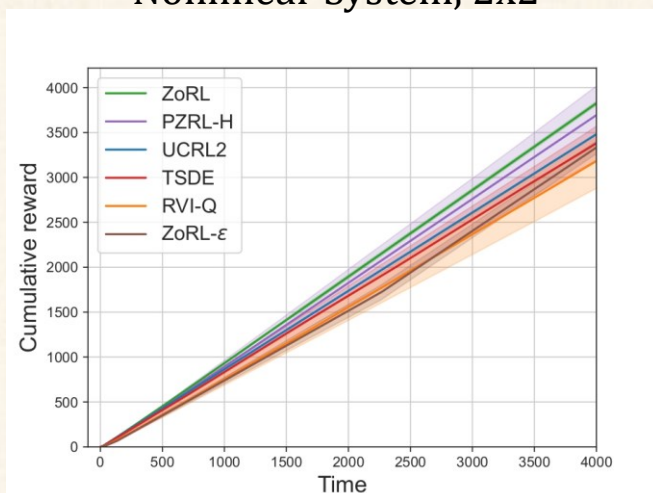
Nonlinear System, 2x2



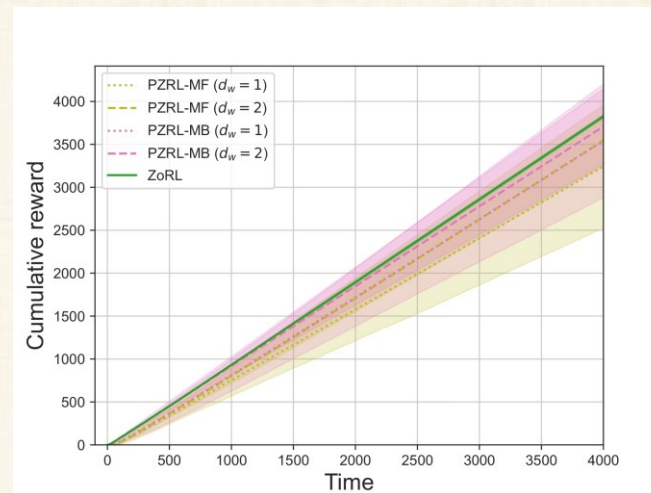
Continuous 2D GridWorld



Truncated Linear System, 2x4



Continuous RiverSwim



Continuous RiverSwim

Zooming in State-Action Space

ZoRL

- I. is statistically efficient
- II. covers a broad class of MDPs
- III. is computationally efficient

Zooming in State-Action Space

ZoRL

- I. is statistically efficient
- II. covers a broad class of MDPs
- III. is computationally efficient

Future directions:

- I. Generalization of zooming idea beyond Lipschitz assumption
- II. Relaxation of other assumptions such as ergodicity
- III. Exploring connection with existing complexity measures

Zooming in State-Action Space

ZoRL

- I. is statistically efficient
- II. covers a broad class of MDPs
- III. is computationally efficient

Future directions:

- I. Generalization of zooming idea beyond Lipschitz assumption
- II. Relaxation of other assumptions such as ergodicity
- III. Exploring connection with existing complexity measures

Preprints:

1. Kar, Avik, and Rahul Singh. "Provably Adaptive Average Reward Reinforcement Learning for Metric Spaces." arXiv preprint arXiv:2410.19919 (2024).
2. Kar, Avik, and Rahul Singh. "Policy Zooming: Adaptive Discretization-based Infinite-Horizon Average-Reward Reinforcement Learning." arXiv preprint arXiv:2405.18793 (2024).

Thank you!