

# Operator World Models for Reinforcement Learning

---

Pietro Novelli

Computational Statistics and Machine Learning  
Istituto Italiano di Tecnologia

# Applications of Reinforcement Learning

- Robotics and Autonomous Systems
- Financial Trading and Portfolio Optimization
- Energy Management and Smart Grids
- Healthcare and Medical Treatment Planning
- Game Strategy and Complex Decision-Making
- Autonomous Vehicles and Transportation

The 2024 Turing Award was assigned yesterday to Andrew Barto and Richard Sutton for developing the conceptual and algorithmic foundations of **reinforcement learning**



# Today's Plan

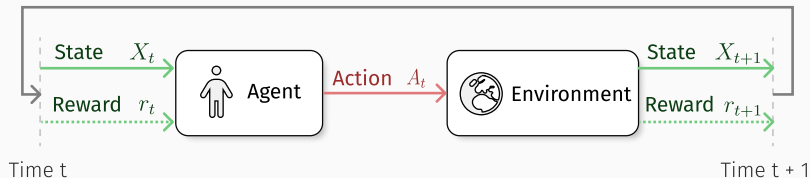
- Describe an original, operator-based, formulation of RL,
- Use it to “solve” RL,
- Introduce approximations to make it practicable,
- Formulate an actual algorithm,
- Watch it in action.

# Problem Setting

---

# Intuition

At time  $t$ , the environment is in the state  $X_t$ . The agent executes an action  $A_t$ . The environment changes its state to  $X_{t+1}$ .



A scalar reward  $r_{t+1} = r(X_t, A_t)$  signal tells us how good/bad the action  $A_t$  was.

# Markov Decision Processes (MDPs)

We consider a **Markov Decision Process (MDP)** characterized by a:

- State space  $\mathcal{X}$
- Action space  $\mathcal{A}$
- A (non-negative) reward  $r : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}_+$
- A transition kernel  $\tau : \Omega = \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{X}) - X_{t+1} \sim \tau(\cdot | x_t, a_t)$ .

**Details:** We assume  $\mathcal{X}$  and  $\mathcal{A}$  Polish spaces and  $\tau, r$  Borel measurable.

$\mathcal{P}(\mathcal{X})$  the space of Borel probability measures on  $\mathcal{X}$ .

**Setting.** We have the freedom to choose what action to take. This is embodied in the notion of a **policy**:

$$\pi : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$$

so that  $\pi(\cdot|x)$  denotes the probability that we will take a specific action when the system is in state  $x \in \mathcal{X}$ .

**Goal.** Find the “best” policy.

# Objective

Given a starting distribution  $\nu \in \mathcal{P}(\mathcal{X})$  and a discount factor<sup>1</sup>  $\gamma \in [0, 1)$ ,

**Goal:** maximize the ( $\gamma$ -discounted) *expected return*

$$J(\pi) = \mathbb{E}_{\nu, \pi, \tau} \left[ \sum_{t=0}^{+\infty} \gamma^t r(X_t, A_t) \right]$$

achieved by *policy*  $\pi$  in the MDP  $(\mathcal{X}, \mathcal{A}, \tau, r)$  starting from  $\nu$ .

**Notation:**

$X_0, A_t$  and  $X_{t+1}$  have laws respectively  $\nu, \pi(\cdot|X_t)$  and  $\tau(\cdot|X_t, A_t)$  for any  $t \in \mathbb{N}$ .

---

<sup>1</sup>This could be generalized, but makes for much more manageable problems.



# Expectations & Objectives

Concretely speaking, we aim to design an algorithm that:

- Returns an **optimal policy** (or a sequence of policies iteratively converging towards the optimum).
- Ideally while establishing some convergence **rates**.

# Reinforcement Learning or Dynamic Programming?

Strictly speaking we talk about

- Dynamic Programming if  $\tau$  is known
- Reinforcement Learning if  $\tau$  is unknown

# Disclaimer

Reinforcement Learning and Dynamic Programming have a relatively long history, dating back to the late 1950s from the work of Bellman<sup>2</sup>, Samuel<sup>3</sup> and Howard<sup>4</sup>.

They are closely related with **Optimal Control** and both very active fields, with a plethora of approaches and techniques developed over the years.

**However**, I am going to blatantly **ignore** all of that and give a very biased and focused talk on a specific and relatively novel perspective on how to tackle them.

---

<sup>2</sup>Bellman, R. *Dynamic Programming*. Princeton University Press. 1957

<sup>3</sup>Samuel, A. L. Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*. 1959

<sup>4</sup>Howard, R. A. *Dynamic Probabilistic Systems*. Wiley. 1960.

## The Operator Way

---

# Markov Decision Processes as Linear Operators

Markov decision processes can be described by **linear operators**<sup>5</sup>.  
How?

As the process is **Markovian**, the next state  $X_{t+1}$  only depends on the present, and not on the past:  $X_{t+1} \sim \tau(\cdot | x_t, a_t)$ .

---

<sup>5</sup>An operator is a map from functions to functions.

# Markov Decision Processes as Linear Operators

Markov decision processes can be described by **linear operators**<sup>5</sup>.  
How?

As the process is **Markovian**, the next state  $X_{t+1}$  only depends on the present, and not on the past:  $X_{t+1} \sim \tau(\cdot | x_t, a_t)$ .

Instead of studying the transition probability  $\tau$ , we can study the integral operator associated to it:

$$(\mathsf{T}f)(x_t, a_t) = \int_{\mathcal{X}} f(x') \tau(dx' | x_t, a_t) = \mathbb{E}[f(X_{t+1}) | x_t, a_t].$$

For any function  $f$  of the state,  $\mathsf{T}f$  is its **expected value in the future** given  $(x_t, a_t)$ .

**T is linear**, since  $\mathsf{T}(f + \alpha g) = \mathsf{T}f + \alpha \mathsf{T}g$ .

---

<sup>5</sup>An operator is a map from functions to functions.

# Why operators?

- They map the problem of studying *non-linear* dynamics on  $\mathcal{X}$  into the problem of studying *linear* operators in function spaces.
- They describe **expected values** in the future, exactly what appears in the RL objective function.
- They can be approximated from data, with proved statistical learning guarantees.

# Operatorial Perspective of Reinforcement Learning

## Transfer Operator.

$\mathbb{T} : B_b(\mathcal{X}) \rightarrow B_b(\Omega)$  such that for any  $f \in B_b(\mathcal{X})$  and  $(x, a) \in \Omega$

$$(\mathbb{T}f)(x, a) = \int_{\mathcal{X}} f(x') \tau(dx'|x, a) = \mathbb{E}[f(X') \mid x, a].$$

## “Policy” Operator.

$P_\pi : B_b(\Omega) \rightarrow B_b(\mathcal{X})$  such that for any  $g \in B_b(\Omega)$  and  $x \in \mathcal{X}$

$$(P_\pi g)(x) = \int_{\mathcal{A}} g(x, a) \pi(da|x) = \mathbb{E}[g(X, A) \mid X = x].$$

## Notation:

$B_b(\mathcal{X})$  space of bounded Borel-measurable functionals on  $\mathcal{X}$ .



## Operatorial Perspective (II)

By applying  $T$  and  $P_\pi$  we have...

Single interaction between  $\pi$  and the MDP starting from  $(x, a)$ :

$$\mathbb{E}[r(X_1, A_1) | X_0 = x, A_0 = a] = (TP_\pi r)(x, a)$$

## Operatorial Perspective (II)

By applying  $T$  and  $P_\pi$  we have...

Single interaction between  $\pi$  and the MDP starting from  $(x, a)$ :

$$\mathbb{E}[r(X_1, A_1) | X_0 = x, A_0 = a] = (TP_\pi r)(x, a)$$

After  $t \in \mathbb{N}$  such steps...

$$\mathbb{E}[r(X_t, A_t) | X_0 = x, A_0 = a] = (TP_\pi)^t r(x, a)$$

## Operatorial Perspective (II)

By applying  $T$  and  $P_\pi$  we have...

Single interaction between  $\pi$  and the MDP starting from  $(x, a)$ :

$$\mathbb{E}[r(X_1, A_1) | X_0 = x, A_0 = a] = (TP_\pi r)(x, a)$$

After  $t \in \mathbb{N}$  such steps...

$$\mathbb{E}[r(X_t, A_t) | X_0 = x, A_0 = a] = (TP_\pi)^t r(x, a)$$

Summing all of them up (with  $\gamma$ -discount):

$$q_\pi(x, a) = \sum_{t=0}^{+\infty} \gamma^t \mathbb{E}[r(X_t, A_t) | x, a] = \sum_{t=0}^{+\infty} (\gamma TP_\pi)^t r = (\text{Id} - \gamma TP_\pi)^{-1} r(x, a)$$

**Why:**  $T$  and  $P_\pi$  are Markov operators, hence their operator norms  $\|T\| = \|P_\pi\| = 1$  and therefore the above Neumann series is convergent.

# Operatorial Perspective (III)

Proposition.<sup>6</sup>

$$J(\pi) = \langle P_\pi q_\pi, \nu \rangle = \langle P_\pi (\text{Id} - \gamma T P_\pi)^{-1} r, \nu \rangle$$

---

<sup>6</sup>With the canonical pairing  $\langle f, \nu \rangle = \int f(x) \nu(dx)$ .

# Operatorial Perspective (III)

Proposition.<sup>6</sup>

$$J(\pi) = \langle P_\pi q_\pi, \nu \rangle = \langle P_\pi (\text{Id} - \gamma TP_\pi)^{-1} r, \nu \rangle$$

**Proof.** Recall the definition of our objective

$$\begin{aligned} J(\pi) &= \mathbb{E}_{\nu, \pi, \tau} \left[ \sum_{t=0}^{+\infty} \gamma^t r(X_t, A_t) \right] \\ &= \mathbb{E}_{\nu, \pi} \left[ \sum_{t=0}^{+\infty} \gamma^t \mathbb{E}[r(X_t, A_t) | X_0, A_0] \right] \\ &= \mathbb{E}_{\nu, \pi} [q_\pi(X_0, A_0)] \\ &= \mathbb{E}_\nu [\mathbb{E}_{\pi(\cdot | X_0)} [q_\pi(X_0, A_0) | X_0]] \\ &= \mathbb{E}_\nu [(P_\pi q_\pi)(X_0)] \\ &= \langle P_\pi q_\pi, \nu \rangle \end{aligned}$$

---

<sup>6</sup>With the canonical pairing  $\langle f, \nu \rangle = \int f(x) \nu(dx)$ .

# Maximizing $J$

Why do we like the operator form for  $J(\pi)$ ? Well, for starters,

$$\max_{\pi} \langle P_{\pi}(\text{Id} - \gamma TP_{\pi})^{-1} r, \nu \rangle$$

is in a much more “standard” form (from an optimization perspective).

Actually... since for any  $\theta \in [0, 1]$  and policies  $\pi_1, \pi_2$ ,

$$P_{\theta\pi_1 + (1-\theta)\pi_2} = \theta P_{\pi_1} + (1 - \theta) P_{\pi_2}$$

The definition of policy operator is **linear** w.r.t. individual policies, so...

...is the problem concave?

Unfortunately, not<sup>7</sup>.

---

<sup>7</sup>Agarwal, Kakade, Lee and Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. JMLR 2021.

<sup>8</sup>We have generalized the definitions of  $J$  and  $q$  to any Markov operator.

# Maximizing $J$

Unfortunately, not<sup>7</sup>.

However, nothing prevents try and minimize it nevertheless!  
In particular, since we can “easily” compute derivatives...

**Lemma.**<sup>8</sup> For any Markov operators  $P, P'$  let  $M = P' - P$ . Then

$$\lim_{h \rightarrow 0} \frac{J(P + hM) - J(P)}{h} = \frac{1}{1 - \gamma} \langle Mq(P), (\text{Id} - \gamma TP)^{-*} \nu \rangle$$

Meaning we could use first-order methods to minimize  $J$ !

---

<sup>7</sup>Agarwal, Kakade, Lee and Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. JMLR 2021.

<sup>8</sup>We have generalized the definitions of  $J$  and  $q$  to any Markov operator.



## Policy Mirror Descent

---

# Mirror Descent<sup>10</sup>

**Mirror Descent (MD).** Optimizing a linear approximation of our objective function, while keeping a small Bregman divergence  $D$  between iterates:

For any  $k \in \mathbb{N}$  the update step is

$$P_{k+1} = \operatorname{argmin}_{P \in \Pi} -\eta \langle \nabla J(P_k), P \rangle + D(P, P_k)$$

If  $J$  were convex, MD guarantees convergence<sup>9</sup>!

(Potential) **Advantage 1:** carefully choosing  $D$  might yield more amenable (e.g. closed-form) solutions for the update step above!

(Potential) **Advantage 2:** As  $J$  depends on  $T$ , for which we will only know an approximation, we will not be “overconfident” in our updates.

---

<sup>9</sup>Under some additional assumption on  $D$  and  $\Pi$

<sup>10</sup>Technically, ascent, but we can just minimize  $-J$

# “Our” Bregman Divergence

So, let's choose our Bregman divergence.

Generalizing (Xiao2022)<sup>11</sup> from the tabular setting<sup>12</sup>, we take

$$D(P_\pi, P_{\pi_k}) = \frac{1}{1 - \gamma} \int \text{KL}(\pi(\cdot|x), \pi_k(\cdot|x)) \rho_{\pi_k}(dx)$$

where:

- $\rho_\pi = (\text{Id} - \gamma \text{TP}_\pi)^{-*} \nu$  is the “occupancy measure” of  $\pi$ ,
- KL is the Kullback-Leibler divergence

---

<sup>11</sup>Xiao. On the convergence rates of policy gradient methods. JMLR 2022.

<sup>12</sup>Tabular setting:  $\mathcal{X}$  and  $\mathcal{A}$  are sets with finite cardinality.

# (Policy) Mirror Descent

**Proposition.** The (Policy) Mirror Descent step can be performed point-wise, namely the iterative sequence of operators  $P_k = P_{\pi_k}$  is such that, for any  $k \in \mathbb{N}$  and any  $x \in \mathcal{X}$

$$\pi_{k+1}(\cdot|x) = \operatorname{argmin}_{p \in \mathbb{P}[\mathcal{A}]} -\eta \langle q_{\pi_k}(\cdot, x), p \rangle + \text{KL}(p, \pi_k(\cdot|x))$$

which has closed-form solution

$$\pi_{k+1}(\cdot|x) = \frac{\pi_k(\cdot|x) e^{\eta q_{\pi_k}(x, \cdot)}}{\int_{\mathcal{A}} \pi_k(a|x) e^{\eta q_{\pi_k}(x, a)} \pi_k(da|x)}$$

## (Policy) Mirror Descent on Finite $\mathcal{A}$

**Assumption:** in the following, we will assume  $\mathcal{A}$  to have finite cardinality.

Then, by applying the PMD update recursively we have

$$\pi_{k+1}(\cdot|x) = \text{SoftMax} \left( \log \pi_0(\cdot|x) + \eta \sum_{j=0}^k q_{\pi_j}(x, \cdot) \right)$$

Where  $\text{SoftMax}(q) = \frac{qe^q}{\sum_{a \in \mathcal{A}} q(x)}$ .

**Note.**

While we could generalize the above to generic  $\mathcal{A}$ , we would be left with an integral at the denominator that we would (likely) be unable to estimate exactly. Studying how such approximation error would propagate will be the subject of future work.

**Theorem<sup>13</sup> (Generalization of (Xiao2022)).** Let  $(\pi_k)_{k \in \mathbb{N}}$  be a sequence generated by PMD with sufficiently large  $\eta > 0$ . Then,

$$\max_{\pi \in \Pi} J(\pi) - J(\pi_k) \leq O(1/k) \quad \forall k \in \mathbb{N}$$

Even if  $J$  is not convex, PMD converges to the global maximum!

---

<sup>13</sup>very informal!

## Towards a practical Algorithm

---

# Taking Stock

**Good news.** We have an algorithm to find the best policy, but...

**Bad news.** For every  $k$  we need to know how to evaluate  $q_{\pi_k}$ . But

$$q_{\pi_k} = (\text{Id} - \gamma \text{TP}_{\pi_k})^{-1} r$$

requires knowledge of the transition operator  $T$ !

## Challenges:

- In Reinforcement Learning (RL) we **do not** know  $\tau$  (or  $T$ )!
- Even in Dynamic Programming, where  $\tau$  is known, it might be too complicated for us to obtain  $T$ !



# Taking Stock

**Good news.** We have an algorithm to find the best policy, but...

**Bad news.** For every  $k$  we need to know how to evaluate  $q_{\pi_k}$ . But

$$q_{\pi_k} = (\text{Id} - \gamma \text{TP}_{\pi_k})^{-1} r$$

requires knowledge of the transition operator  $T$ !

## Challenges:

- In Reinforcement Learning (RL) we **do not** know  $\tau$  (or  $T$ )!
- Even in Dynamic Programming, where  $\tau$  is known, it might be too complicated for us to obtain  $T$ !

**Idea:** let's approximate  $q_{\pi_k}$  with some  $\hat{q}_{\pi_k}$  that is more amenable to practical manipulations!

# Covergence of “Approximate” PMD

**Theorem<sup>14</sup> (Generalization of (Xiao2022)).** Let  $(\pi_k)_{k \in \mathbb{N}}$  be a sequence generated by the “approximate” PMD step

$$\pi_k(\cdot|x) = \text{SoftMax} \left( \log \pi_0(\cdot|x) + \sum_{j=0}^k \hat{q}_{\pi_k}(x, \cdot) \right)$$

where  $\hat{q}_{\pi_k}$  are such that  $\|\hat{q}_{\pi_k} - q_{\pi_k}\|_{\infty} \leq \epsilon_k$  for some  $\epsilon_k > 0$ . Then,

$$\max_{\pi \in \Pi} J(\pi) - J(\pi_k) \leq O \left( \frac{1 + \sum_{j=0}^k \epsilon_j}{k} \right) \quad \forall k \in \mathbb{N}$$

If we can control the  $\epsilon_j$  (e.g. such that  $\epsilon = O(1/j)$ ), then “approximate” PMD converges to the global maximum!

---

<sup>14</sup>Again, very informal!

## Approximating $q_\pi$ with World Models

---

# Approximating $q_\pi$ using World Models

The operator perspective on  $q_\pi$  offers a direct strategy to define a  $\hat{q}_\pi$

$$\hat{q}_\pi = (\text{Id} - \gamma \hat{T} P_\pi)^{-1} \hat{r}$$

In other words, we need to approximate (or learn!):

- The one-step update of the environment — a “world model”  $\hat{T}$ .
- The immediate reward function  $\hat{r}$ .

**Note.** We have also to ensure that the definition of  $\hat{q}_\pi$  makes sense...

# Reproducing Kernel Hilbert Spaces

We will rely on standard machine learning tools: [kernel methods](#).

**Reward Function.** Let  $\psi : \Omega \rightarrow \mathcal{G}$  be a feature map of a reproducing kernel Hilbert space<sup>15</sup> (rkhs)  $\mathcal{G}$ . Then, given  $n \in \mathbb{N}$  points  $(x_i, a_i)_{i=1}^n$ ,

$$r_n = \operatorname{argmin}_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \left( \langle g, \psi(x_i, a_i) \rangle_{\mathcal{G}} - r(x_i, a_i) \right)^2 + \lambda \|g\|_{\mathcal{G}}^2$$

where  $\lambda > 0$  is a regularization parameter.

**Notation.** We will replace  $\hat{r}$  with  $r_n$  to highlight the dependency on  $n$ .

---

<sup>15</sup>Namely,  $\mathcal{G}$  is a space of functions  $g(x, a) = \langle g, \psi(x, a) \rangle$

# Learning the Reward Function

**Ridge Regression.** The quantity  $r_n$  admits closed-form solution

$$r_n = S_n^* b = \sum_{i=1}^n b_i \psi(x_i, a_i) \quad \text{where} \quad b = (K + \lambda \text{Id})^{-1} y$$

where

- $y \in \mathbb{R}^n$  is the vector with entries  $y_i = r(x_i, a_i)$ ,
- $K \in \mathbb{R}^{n \times n}$  the “kernel matrix” with entries

$$K_{ij} = k((x_i, a_i), (x_j, a_j)) = \langle \psi(x_i, a_i), \psi(x_j, a_j) \rangle_{\mathcal{G}}$$

- $S_n : \mathcal{G} \rightarrow \mathbb{R}^n$  such that  $S_n : g \mapsto (g(x_i, a_i))_{i=1}^n$

**Take-home message.** We have a “finite” representation of  $r_n$  that fits into a machine that we can use in practice!

# Conditional Mean Embeddings

Can we do the same thing for  $T$ ? Yes, if we restrict to suitable spaces...

**Remark.** Let  $\mathcal{G}$  and  $\mathcal{F}$  two rkhs over  $\Omega$  and  $\mathcal{X}$  with feature maps

$\psi : \Omega \rightarrow \mathcal{G}$  and  $\varphi : \mathcal{X} \rightarrow \mathcal{F}$  respectively.

Then, if the restriction of  $T$  to  $\mathcal{F}$  takes values in  $\mathcal{G}$ .

$$(T|_{\mathcal{F}})^* \psi(x, a) = \int \varphi(x') \tau(dx'|x, a) \quad \forall (x, a) \in \Omega$$

In other words,  $\psi(x, a)$  is mapped to the conditional expectation of  $\psi(x')$ !

**Def.**  $T|_{\mathcal{F}}$  is known as the conditional mean embedding (CME) of  $\tau$ .

**Idea.** If we can sample from  $\tau$ , we can collect a dataset  $(x_i, a_i, \varphi(x'_i))_{i=1}^n$  and learn  $T|_{\mathcal{F}}$  like we did for  $r_n$ .

# Conditional Mean Embedding

We formulate the learning problem over Hilbert-Schmidt operators,

$$\tilde{T}_n = \operatorname{argmin}_{W \in \text{HS}(\mathcal{F}, \mathcal{G})} \frac{1}{n} \sum_{i=1}^n \|W^* \psi(x_i, a_i) - \varphi(x'_i)\|_{\mathcal{F}}^2 + \lambda \|W\|_{\text{HS}}^2$$

which yields the closed-form solution

$$\tilde{T}_n = S_n^*(K + \lambda \text{Id})^{-1} Z_n$$

with  $Z_n : \mathcal{F} \rightarrow \mathbb{R}^n$  such that  $Z_n : f \mapsto (f(x'_i))_{i=1}^n$ .

**Normalization.** We then take  $T_n = \frac{\tilde{T}_n}{\|\tilde{T}_n\|}$  to ensure  $\|T_n\| = 1$

**Take-home message 2.** We have a “finite” representation of  $T_n$  that fits into a machine that we can use in practice!



POWR

---

# POWR Algorithm

We have **Policy Mirror Descent with Operator World-models for RL**.

- Collect a dataset  $(x_i, a_i, x'_i)_{i=1}^n$  of sample transitions to learn  $T_n$  (analogously for  $r_n$ ).
- Choose  $\pi_0$ , for example  $\pi_0(\cdot|x)$  uniform for any  $x \in \mathcal{X}$ .
- For  $k = 0, \dots$ ,
  - Let  $q_{\pi_k, n} = (\text{Id} - \gamma T_n P_{\pi_k})^{-1} r_n$
  - Let  $\pi_{k+1} = \text{SoftMax} \left( \log \pi_0 + \sum_{j=0}^k q_{\pi_j, n} \right)$
- **Return:**  $\pi_k$  for any  $k \in \mathbb{N}$

---

**Algorithm 1** POWR: POLICY MIRROR DESCENT WITH OPERATOR WORLD-MODELS FOR RL
 

---

**Input:** Dataset  $(x_i, a_i, x'_i, r_i)_{i=1}^n$ , discount factor  $\gamma \in (0, 1)$ , step size  $\eta > 0$ , kernel function  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  with  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  as in Proposition 4, initial weights  $C_0 = 0 \in \mathbb{R}^{n \times |\mathcal{A}|}$ .

*/\* World Model Learning \*/*  
 let  $E \in \mathbb{R}^{n \times |\mathcal{A}|}$  with rows  $E_i = \text{ONEHOT}_{|\mathcal{A}|}(a_i)$ .  
 let  $K_\lambda \in \mathbb{R}^{n \times n}$  such that  $K_{ij} = k(x_i, x_j)\delta_{a_i=a_j} + n\lambda\delta_{ij}$   
 let  $H \in \mathbb{R}^{n \times n}$  such that  $H_{ij} = k(x'_i, x_j)$   
 compute  $K_\lambda^{-1}$  and  $b = K_\lambda^{-1}y$  with  $y = (r_1, \dots, r_n) \in \mathbb{R}^n$

*/\* Policy Mirror Descent \*/*  
 for  $t = 0, 1, \dots, T-1$  do:  
      $\pi_{t+1} = \text{SOFTMAX}(\eta H C_t) \in \mathbb{R}^{n \times |\mathcal{A}|}$   
      $M_{\pi_{t+1}} = H \odot (\pi_{t+1} E^\top) \in \mathbb{R}^{n \times n}$   
      $C_{t+1} = C_t + \text{diag}(c)E$  with  $c = (\text{Id} - \gamma K_\lambda^{-1} M_{\pi_{t+1}})^{-1}b$   
 end for

**return**  $\pi_T: \mathcal{X} \rightarrow \Delta(\mathcal{A})$  such that  $\pi_T(x) = \text{SOFTMAX}(\eta H_x C_T)$  with  $H_x = (k(x, x_i))_{i=1}^n \in \mathbb{R}^n$ .

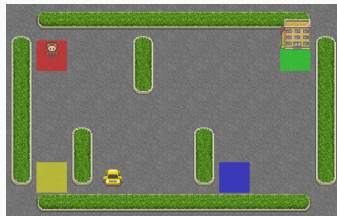
---

# Experiments

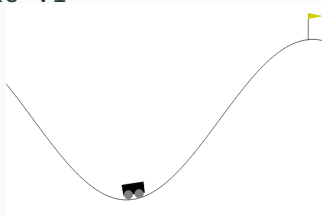
For now, we have tried POWR on very small-scale/toy environments.



Frozen Lake-v1

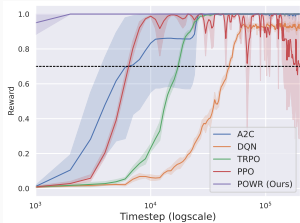


Taxi-v3

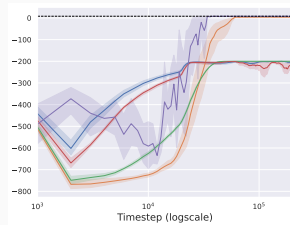


Mountain Car-v0

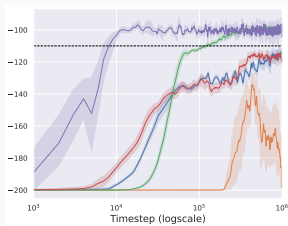
# Empirical POWR Sample Efficiency



Frozen Lake-v1



Taxi-v3



Mountain Car-v0

# Recap

We...

- Set off to tackle sequential decision making problems...
- Saw how the operator-based perspective offered interesting insights.
- Saw that policy mirror descent converges to the global maximum **provided we can approximate  $q_\pi$** .
- Proposed an estimator for  $q_\pi$  in terms of a “world model”  $T_n$  (and an estimate for the reward  $r_n$ ).
- Observed that POWR actually works well in practice.

# Open Questions

- (Scaling up) How well does POWR work on more challenging environments?
- (Efficiency) The usual suspects, Nystrom, Random Features, etc.
- (Infinite Actions) Can we adapt POWR to infinite action spaces?
- (Exploration Vs Exploitation) How to choose the distribution  $\rho$  from which we obtain the dataset to train  $T_n$  and  $r_n$ ?

Questions?



## EXTRA: Theoretical analysis of POWR

---

# Does POWR “Work”?

Two main questions:

- Are POWR’s iterates  $(\mathcal{G}, \mathcal{F})$ -compatible? And why do we care?
- When (if ever) does POWR converge?

# Guaranteeing Compatibility - Factoring $\mathcal{F}$ and $\mathcal{G}$

We restrict to the following choice for  $\mathcal{F}$  and  $\mathcal{G}$ :

- $\mathcal{H}$  be a rkhs with feature map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$ .
- $\mathcal{F} = \mathcal{H} \otimes \mathcal{H}$  with  $\varphi(x) = \phi(x) \otimes \phi(x)$ ,
- $\mathcal{G} = \mathbb{R}^{|\mathcal{A}|} \otimes \mathcal{H}$  with<sup>16</sup>  $\psi(x, a) = e_a \otimes \phi(x)$ .

Then – recalling that  $\mathcal{A}$  is a finite set – we have the following,

**Proposition.** A policy  $\pi$  is  $(\mathcal{G}, \mathcal{F})$ -compatible if and only if there exist  $p_a \in \mathcal{H}$  such that  $\pi(a|\cdot) = \langle p_a, \phi(\cdot) \rangle_{\mathcal{H}}$  for and  $a \in \mathcal{A}$ .

It is enough to check that all  $\pi(a|\cdot)$  “belong” to  $\mathcal{H}$  to guarantee  $(\mathcal{F}, \mathcal{G})$ -compatibility!

---

<sup>16</sup>Here,  $e_a$  is the  $a$ -th element of the canonical basis of  $\mathbb{R}^{|\mathcal{A}|}$  (assume an order on  $\mathcal{A}$ ).

# Sobolev Spaces to the Rescue!

**Theorem.** Let  $\mathcal{X} \subset \mathbb{R}^d$  be compact,  $\mathcal{H} = W^{2,s}(\mathcal{X})$  the Sobolev space with smoothness  $s > d/2$ . Let  $\pi_0(a|\cdot) \propto e^{\eta q_0(\cdot, a)}$  for some  $q_0(\cdot, a) \in \mathcal{H}$  for all  $a \in \mathcal{A}$ .

$\implies$  all iterates produced by POWR are  $(\mathcal{F}, \mathcal{G})$ -compatible.

**Proof sketch.** The key is to show recursively that

- If  $\pi_k$  is  $(\mathcal{G}, \mathcal{F})$ -compatible, then the approximate  $q_{\pi_k, n}$  belong to  $\mathcal{H}$  and,
- The SoftMax operator applied to previous  $q_{\pi_k, n}$  yields a  $(\mathcal{G}, \mathcal{F})$ -compatible policy  $\pi_{k+1}$

# Does POWR “Work”?

Two main questions:

- Are POWR’s iterates  $(\mathcal{G}, \mathcal{F})$ -compatible? **Yes!**
- When (if ever) does POWR converge?

# Convergence of POWR - Assumptions

POWR converges under suitable regularity assumptions...

**Assumption (Strong Source Condition).** There exists <sup>17</sup>  $\rho \in \mathcal{P}(\Omega)$  s.t.

$$\|(T|_{\mathcal{F}})^* C_{\rho}^{-\beta}\|_{\text{HS}} < +\infty \quad \text{and} \quad \|C_{\rho}^{-\beta} r\|_{\mathcal{G}} < +\infty,$$

for some  $\beta > 0$ , where  $C_{\rho} = \sum_{a \in \mathcal{A}} \int_{\mathcal{X}} \psi(x, a) \otimes \psi(x, a) \rho(dx, a)$ .

## Notes.

- This is a stronger version of the standard assumption used in supervised learning settings.
- We need it because we will  $T_n \rightarrow T$  and  $r_n \rightarrow r$  to be in a stronger norm than usual.

---

<sup>17</sup>We are implicitly asking  $r \in \text{range}(C_{\rho}^{\beta})$  and  $\text{range}(T|_{\mathcal{F}}) \subseteq \text{range}(C_{\rho}^{\beta})$ .

# Convergence of POWR

**Theorem.** Let  $\rho$  satisfy the Strong Source Condition. Let the world-model  $T_n$  and reward  $r_n$  estimators learned from a dataset  $(x_i, a_i, x'_i)_{i=1}^n$  where  $(x_i, a_i)$  are independently sampled from  $\rho$  and  $x'_i \sim \tau(\cdot | x_i, a_i)$  for  $i = 1, \dots, n$ . Then, for any  $\delta \in (0, 1)$ , the iterates produced by POWR converge to the optimal return as

$$\max_{\pi \in \Pi} J(\pi) - J(\pi_k) \leq O\left(\frac{1}{K} + \delta n^{-\frac{\beta}{2+2\beta}}\right)$$

with probability not smaller than  $1 - 4e^{-\delta}$

- **Good news:** it converges!
- **Bad news:** maybe not that fast...

# Proof Sketch

## Proof sketch.

- We know already that PMD with approximate  $q_{\pi,n}$  converges with rate  $O(1/k + \epsilon)$ , if  $\|q_{\pi_k,n} - q_{\pi}\| \leq \epsilon$  uniformly wrt  $k \in \mathbb{N}$ .
- The following Lemma gives us an idea of how to control  $\epsilon$ :  
**Lemma.** Assume  $T|_{\mathcal{F}} : \mathcal{F} \rightarrow \mathcal{G}$ . Then

$$\|q_{\pi,n} - q_{\pi}\|_{\infty} \leq O(\|r_n - r\|_{\infty} + \|r\|_{\infty} \|T_n - T|_{\mathcal{F}}\|_{\text{HS}})$$

- Bounding bounding  $\epsilon$  boils down to controlling the approximation error of  $r_n$  and  $T_n$  in  $\|\cdot\|_{\infty}$  norm. This is a supervised setting and we can therefore borrow refined results from the literature<sup>18</sup>

---

<sup>18</sup>for example Fischer and Steinwart. Sobolev norm learning rates for regularized least-squares algorithms. JMLR 2020.



I am hiding a lot of details/questions:

- Constants depending on the key quantities of the problem.
- Minimum sample size  $n$  required to make everything work.
- How to choose the step size  $\eta$ ?
- How to choose  $\rho$ ?
- ...