Policy gradient methods for RL in general spaces

David Šiška

Tutorial for Bridging Stochastic Control And Reinforcement Learning

7th November 2025



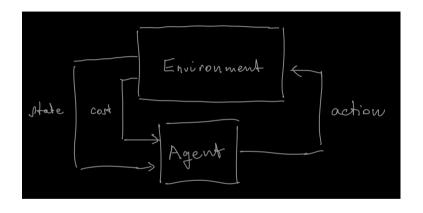


Slides Notes

Outline

- Reinforcement Learning (RL) and its (entropy regularized) MDP formulation
 - Relative entropy and its key properties
 - Bellman principle for (entropy regularized) MDPs
- Classical Policy Gradient (PG)
 - Performance difference lemma and policy gradient theorem
 - Difficulty of convergence analysis due to lack of convexity
 - Polyak-Lojasiewicz (PL) gradient dominance condition
- Mirror descent
 - Role of Performance difference as convexity and L-smoothness
 - Convergence rate MDP case with inexact advantage
- PPO algorithm
 - Convergence of FR-PPO variant

Reinforcement Learning (RL)

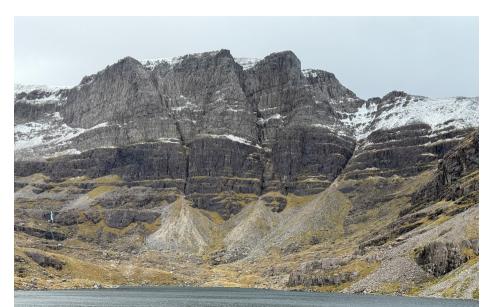


RL Aim: learn to interact with an environment in an optimal (cost minimizing) way.

Data: $(s_t, a_t, c_t, s_{t+1}, a_{t+1}, \ldots)$.

Mathematical abstraction: MDP.

MDPs and relaxed, regularized MDPs



Key MDP / RL results

Overview of RL [Sutton and Barto, 2018] and results in discrete state-action spaces

- Classical policy gradient [Sutton et al., 1999].
- Natural policy gradient [Kakade, 2001].
- Actor-critic method [Haarnoja et al., 2018].
- Mirror descent method [Tomar et al., 2020].
- Convergence of classical PG in tabular setting [Mei et al., 2021].

Continuous state-action spaces: [Doya, 2000], [Van Hasselt, 2012], [Manna et al., 2022].

Entropy regularised: [Haarnoja et al., 2017, Geist et al., 2019].

Infinite-horizon Markov decision problem (S, A, P, c, γ) :

- *S* is the state space, *A* is the action space
- ullet $P \in \mathcal{P}(S|S \times A)$ is the transition probability kernel
- $c \in B_b(S \times A)$ is a cost function, and γ discount factor
- $H_n := (S \times A)^n \times S$ is the space of admissible histories

Aim: minimise the objective over policies $\alpha = (\alpha_n)_{n \in \mathbb{N}}$ s.t. $\alpha_n : H_n \to A$ measurable:

$$V^{\alpha}(s) = \mathbb{E}_{s}^{\alpha} \sum_{n=0}^{\infty} \gamma^{n} c(s_{n}, a_{n}), \qquad (1)$$

with $a_n := \alpha_n(h_n)$, $h_n = (s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n)$ and $s_{n+1} \sim P(\cdot | s_n, a_n)$, $s_0 = s$.

Relaxed and regularized formulation of the MDP

Infinite-horizon Markov decision model (S, A, P, c, γ) :

- S is the state space, A is the action space,
- $P \in \mathcal{P}(S|S \times A)$ is the transition probability kernel,
- $c \in B_b(S \times A)$ is a cost function, and γ a discount factor,
- $H_n := (S \times A)^n \times S$ is the space of admissible histories,
- $\tau \geq 0$ strenght of entropy regularizer,
- for $\mu', \mu \in \mathcal{P}(A)$ define $\mathsf{KL}(\mu'|\mu) = \int_A \mathsf{In} \, \frac{\mathrm{d}\mu'}{\mathrm{d}\mu}(a)\mu'(da)$ if $\mu' \ll \mu$, and $+\infty$ otherwise.

Aim: minimise over relaxed policies $\pi = (\pi_n)_{n \in \mathbb{N}}$ s.t. $\pi_n : H_n \to P(A)$ measurable the objective:

$$V_{\tau}^{\pi}(s) = \mathbb{E}_{s}^{\pi} \left[\sum_{n=0}^{\infty} \gamma^{n} \left(\int_{A} c(s_{n}, \mathbf{a}) \, \pi_{n}(d\mathbf{a}) + \tau \, \mathsf{KL}(\pi_{n} | \mu) \right) \right] \in \mathbb{R} \cup \{+\infty\} \,, \tag{2}$$

with $\pi_n := \pi_n(h_n)$, $h_n = (s_0, a_0, \dots, s_{n-1}, a_{n-1}, s_n)$ and $s_{n+1} \sim \int_A P(\cdot|s_n, a) \pi_n(da)$, $s_0 = s$.

Kullback-Leibler divergence aka relative entropy



Relative entropy - definition and basics

If $\nu, \mu \in \mathcal{P}(A)$ and if $\mu(B) = 0 \implies \nu(B) = 0$ for every $B \in \mathcal{B}(A)$ then we say ν is absolutely continuous w.r.t. μ (notation $\nu << \mu$).

For $\mu \in \mathcal{P}(A)$ define

$$\mathcal{P}(A) \ni \nu \mapsto \mathsf{KL}(\nu|\mu) = \begin{cases} \int_A \ln rac{\mathrm{d}\nu}{\mathrm{d}\mu} \, \nu(\mathit{da}) & \text{if } \nu << \mu \,, \\ +\infty & \text{otherwise} \,. \end{cases}$$

Note that

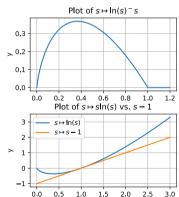
$$\int_{A} \left(\ln rac{\mathrm{d}
u}{\mathrm{d} \mu}
ight)^{-}
u(extit{da}) = \int_{A} \left(\ln rac{\mathrm{d}
u}{\mathrm{d} \mu}
ight)^{-} rac{\mathrm{d}
u}{\mathrm{d} \mu} \, \mu(extit{da})$$

and $s\mapsto (\ln s)^-s\geq 0$ is bounded for $s\geq 0$, so KL is well defined.

Moreover $s \ln s \geq s-1$ for $s \geq 0$ (with equality only if s=1) and so

$$\mathsf{KL}(\nu|\mu) = \int_A \left(\, \mathsf{In} \, \tfrac{\mathrm{d}\nu}{\mathrm{d}\mu} \right) \tfrac{\mathrm{d}\nu}{\mathrm{d}\mu} \, \mu(\textit{da}) \geq \int_A \left(\tfrac{\mathrm{d}\nu}{\mathrm{d}\mu} - 1 \right) \mu(\textit{da}) = 0 \, ,$$

with equality only if $\frac{\mathrm{d}\nu}{\mathrm{d}\mu}=1$ i.e. if $\nu=\mu$.



Relative entropy - variational formula

Useful identity

$$\mathsf{KL}(\nu|\mu) - \mathsf{KL}(\nu'|\mu) = \mathsf{KL}(\nu|\nu') + \int_{A} \ln \frac{\mathrm{d}\nu'}{\mathrm{d}\mu}(a)(\nu - \nu')(da). \tag{3}$$

which holds for any $\nu, \nu' \in \mathcal{P}(A)$ for which the quantities in the identity are finite.

Variational formula: for $f \in B_b(A)$:

$$\inf_{
u \in \mathcal{P}(\mathcal{A})} \left(\int_{\mathcal{A}} f \, d
u + \mathsf{KL}(
u | \mu) \right) = - \ln \int_{\mathcal{A}} e^{-f} \mu(d\mathsf{a}) \, ,$$

and if

$$\frac{\mathrm{d}\nu^*}{\mathrm{d}\mu}(a) = \frac{\mathrm{e}^{f(a)}}{\int_A \mathrm{e}^{-f(a')}\mu(da')}$$

then $u^* = \operatorname{argmin}_{\nu \in \mathcal{P}(A)} \left(\int_A f \, d\nu + \operatorname{KL}(\nu | \mu) \right).$

Relative entropy - dual formulation and convexity

Donsker-Varadhan variational formula

$$\mathsf{KL}(
u|\mu) = \sup_{g \in C_b(A)} \left(\int_A g(a) \, \nu(da) - \ln \int_A e^{g(a)} \, \mu(da) \right)$$

and

$$\mathsf{KL}(
u|\mu) = \sup_{\psi \in B_b(A)} \left(\int_A \psi(\mathsf{a}) \,
u(\mathsf{d}\mathsf{a}) - \mathsf{ln} \int_A e^{\psi(\mathsf{a})} \, \mu(\mathsf{d}\mathsf{a}) \right) \, .$$

N.B. for fixed g

$$(
u,\mu)\mapsto \int_A g(a)\,
u(da)-\ln\int_A \mathrm{e}^{g(a)}\,\mu(da)$$

is convex. As a supremum over such g

• $\mathcal{P}(A) \times \mathcal{P}(A) \ni (\nu, \mu) \mapsto \mathsf{KL}(\nu|\mu)$ is convex, lower-semicontinuous.

Moreover

• For fixed $\mu \in \mathcal{P}(A)$ we have

$$\{\nu \in \mathcal{P}(A) : \mathsf{KL}(\nu|\mu) < \infty\} \ni \nu \mapsto \mathsf{KL}(\nu|\mu)$$

strictly convex, from strict convexity of $[0,\infty) \ni s \mapsto s \ln s \in \mathbb{R}$.

All from [Dupuis and Ellis, 1997, Ch. 1, Sec. 4].

Bellman principle aka Dynamic programming principle (DPP)



DPP with $\tau > 0$

Recall $H_n := (S \times A)^n \times S$ is the space of admissible histories.

Let $V_{ au}^*: \mathcal{S} o \mathbb{R}$ be

$$V_{\tau}^{*}(s) = \inf_{\pi} V_{\tau}^{\pi}(s), \quad \forall s \in S,$$

$$\tag{4}$$

where infimum is over policies $\pi = (\pi_n)_{n \in \mathbb{N}}$ s.t. $\pi_n : H_n \to P(A)$ is measurable.

Theorem 1 (Dynamic programming principle)

Let $\tau > 0$. The optimal value function V_{τ}^* is the unique bounded solution of

$$V_{ au}^*(s) = \inf_{m \in \mathcal{P}(A)} \int_A \left(c(s,a) + au \ln rac{\mathrm{d} m}{\mathrm{d} \mu}(a) + \gamma \int_S V_{ au}^*(s') P(ds'|s,a)
ight) m(da), \quad orall s \in S$$

DPP consequences for $\tau > 0$

For all $s \in S$,

$$V_{ au}^*(s) = - au \ln \int_A \exp\left(-rac{1}{ au}Q_{ au}^*(s,a)
ight) \mu(da),$$

where $Q^* \in B_b(S \times A)$ is defined by

$$Q_{ au}^*(s,a) = c(s,a) + \gamma \int_{\mathcal{S}} V_{ au}^*(s') P(ds'|s,a) \,, \quad orall (s,a) \in \mathcal{S} imes \mathcal{A}.$$

Moreover, there is an optimal policy $\pi_{\tau}^* \in \mathcal{P}_{\mu}(A|S)$ given by

$$\pi_{\tau}^*(da|s) = \exp\left(-(Q_{\tau}^*(s,a) - V_{\tau}^*(s))/\tau\right)\mu(da), \quad \forall s \in S.$$

Let

$$\Pi_{\mu} = \left\{\pi \in \mathcal{P}(A|S) : \ln rac{d\pi}{dx} \in B_b(S \times A) \right\}.$$

Then

$$\inf_{\pi} V^{\pi}_{ au} = V^*_{ au}(s) = \inf_{\pi \in \Pi} V^{\pi}_{ au}$$
 .

Finally, for each $\pi \in \Pi_{\mu}$, we define the *Q*-function $Q_{\tau}^{\pi} \in B_b(S \times A)$ by

$$Q^\pi_ au(s,a) = c(s,a) + \gamma \int_{\mathbb{R}} V^\pi_ au(s') P(ds'|s,a) \,.$$

(6)

(5)

DPP consequences for $\tau > 0$

Proposition 1

Let $f \in B_b(S \times A)$ and $\pi \in \Pi_\mu$ be such that $\pi(da|s) = \frac{\exp(f(s,a))\mu(da)}{\int_A \exp(f(s,a'))\mu(da')}$ for all $s \in S$. Then

$$\left\| \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu} \right\|_{B_b(S \times A)} \leq 2\|f\|_{B_b(S \times A)} \,, \quad \|V_\tau^\pi\|_{B_b(S)} \leq \frac{1}{1-\gamma} \left(\|c\|_{B_b(S \times A)} + 2\tau \|f\|_{B_b(S \times A)} \right) \,.$$

Proof. As $\mu(A) = 1$, for all $g \in B_b(S \times A)$ and $s \in S$,

$$\begin{split} & \ln \int_A \exp(g(s,a')) \mu(da') \leq \ln \left(e^{\|g\|} B_b(S \times A) \, \mu(A)\right) = \|g\|_{B_b(S \times A)} \,, \\ & \ln \int_s \exp(g(s,a')) \mu(da') \geq \ln \left(e^{-\|g\|} B_b(S \times A) \, \mu(A)\right) = -\|g\|_{B_b(S \times A)} \,. \end{split}$$

Then, for all $(s, a) \in S \times A$, using $\ln \frac{d\pi}{d\mu}(a|s) = f(s, a) - \ln \int_A \exp(f(s, a'))\mu(da')$,

$$\left|\ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a|s)\right| \leq |f(s,a)| + \left|\ln \int_{\mathbb{R}} \exp(f(s,a'))\mu(da')\right| \leq 2\|f\|_{B_h(S\times A)},$$

which implies that

$$\left|\mathbb{E}_s^{\pi}\left[\sum_{t=0}^{\infty}\gamma^t\bigg(\tau\ln\frac{\mathrm{d}\pi}{\mathrm{d}\mu}\big(a_t|s_t\big)\bigg)\right]\right|\leq 2\tau\|f\|_{\mathcal{B}_b(S\times A)}\sum_{t=0}^{\infty}\gamma^t=\frac{2\tau\|f\|_{\mathcal{B}_b(S\times A)}}{1-\gamma}\;.$$

The rest follows as usual.

Lemma 2

Let $\tau>0$ and $\pi\in\Pi_{\mu}$. The value function V^{π}_{τ} is the unique bounded solution of the on-policy Bellman equation:

$$V^\pi_ au(s) = \int_{\mathbb{R}} \Big(c(s,a) + au \ln rac{\mathrm{d}\pi}{\mathrm{d}\mu}(a|s) + \gamma \int_{\mathbb{R}} V^\pi_ au(s') P(ds'|s,a) \Big) \pi(da|s), \quad orall s \in S \,.$$

Note that from this and defn. of the Q-function (6) we have for all $\pi \in \Pi_{\mu}$ and $s \in S$ that

$$V_{\tau}^{\pi}(s') = \int_{A} \left(Q_{\tau}^{\pi}(s', a') + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a'|s') \right) \pi(da'|s'), \quad \forall s \in S.$$
 (7)

Using this in the defn. of the Q-function (6) we have the on policy Q-Bellman equation

$$Q_{\tau}^{\pi}(s,a) = c(s,a) + \gamma \int_{S} \int_{A} \left(Q_{\tau}^{\pi}(s',a') + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a'|s') \right) \pi(da'|s') P(ds'|s,a) , \ \forall (s,a) . \tag{8}$$

DPP with $\tau = 0$

Assumption 2

- The kernel $P \in \mathcal{P}(S|S \times A)$ is strongly continuous, that is: for every $v \in B_b(S)$ (bounded and measurable) the function $w(s, a) = \int_S v(s')P(\mathrm{d}s'|s, a)$ is bounded and measurable as a function from $S \times A$ to \mathbb{R} .
- ② The cost function $c \in B_b(S \times A)$ is lower semi-continuous and inf-compact on $S \times A$ i.e. for any $s \in S$ and any $l \in \mathbb{R}$ the set $\{a \in A : c(s, a) \leq l\}$ is compact.

Theorem 3 (Dynamic programming principle, $\tau = 0$)

Let Assumption 2 hold. Then the optimal value function $V^* \in B_b(S)$ is the unique solution of the Bellman equation

$$V^*(s) = \min_{a \in A} \left[c(s, a) + \gamma \int_{S} V^*(s') P(\mathrm{d}s'|s, a) \right]. \tag{9}$$

Moreover, writing $Q^*(s,a) = c(s,a) + \gamma \int_S V^*(s') P(\mathrm{d}s'|s,a)$, there exists a measurable function $f^*: S \to A$ called a selector such that $f^*(s) \in \operatorname{argmin}_{a \in A} Q^*(s,a)$ and the induced policy $\pi^* \in \mathcal{P}(A|S)$ defined by $\pi^*(\mathrm{d}a|s) = \delta_{f^*(s)}(\mathrm{d}a)$ for all $s \in S$ satisfies $V^* = V^{\pi^*}$.

Proof. [Hernández-Lerma and Lasserre, 2012, Theorem 4.2.3].

Proposition 3

Let $\pi(da|s) \in \mathcal{P}(A|S)$. Then

$$\|V_0^{\pi}\|_{B_b(S)} \leq \frac{\|c\|_{B_b(S\times A)}}{1-\gamma}$$
.

Proof. Exercise, start with (1) which is definition of V_0^{π} .

Lemma 4

Let $\pi \in \mathcal{P}(A|S)$. The value function V_0^{π} is the unique bounded solution of the on-policy Bellman equation:

$$V_0^\pi(s) = \int_{\mathbb{A}} \left(c(s,a) + \gamma \int_{\mathbb{C}} V_0^\pi(s') P(ds'|s,a) \right) \pi(da|s), \quad \forall s \in S.$$

What does "solving our RL problem" mean?

We will say we've "solved our RL problem" if we can find a near optimal policy for the MDP under the assumptions that:

- We do **not** have access to costs c and transitions $P \in \mathcal{P}(S|S \times A)$.
- We choose $\gamma > 0$, $\tau \geq 0$
- We have access to a simulator of the environment and we can repeatedly use it cost-free.
- The simulator will initialise at $s \sim \rho \in \mathcal{P}(S)$ of its choice and will run until termination or until we reset it.

Key meta-algorithms

Policy gradient (PG)

- 1: Initialize environment, parametrized policy π_{θ} ,
- 2: for $n = 0, 1, \dots, N_{\text{episodes}}$ do
- 3: Clear memory buffer
- 4: **for** $t = 0, 1, ..., N_{\text{steps in episode}}$ **do**
- 5: Observe state s_t
 - Sample action $a_t \sim \pi_{\theta_n}(a_t|s_t)$
- 7: Execute a_t in environment, accept cost c_t , new state s_{t+1}
- 8: Store $(s_t, a_t, c_t, \log \pi_{\theta_n}(a_t|s_t), V(s_t))$
- 9: $t \leftarrow t + 1$, $s_t \leftarrow s_{t+1}$ in memory.
- 10: end for
- 11: Estimate $\nabla_{\theta} V^{\pi_{\theta}}$ from memory data
- 12: Update policy parameters

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} V^{\pi_{\theta}}$$
.

13: end for

Q-learning

- 1: Initialize environment, parametrized state-action value Q_{θ} ,
- 2: for $n = 0, 1, \ldots, N_{\text{episodes}}$ do
- 3: Make space in memory buffer
- 4: for $t = 0, 1, ..., N_{\text{steps in episode}}$ do
- : Observe state s_t
- 6: Take $a_t \in \operatorname{argmin}_{a \in A} Q_{\theta}(s_t|a)$
- 7: Execute a_t in environment, accept cost c_t , new state s_{t+1}
- 8: Store (s_t, a_t, c_t, s_{t+1}) in memory
- 9: $t \leftarrow t+1$, $s_t \leftarrow s_{t+1}$
- 10: end for
- 11: Sample $(s_j, a_j, c_j, s_{j+1})_{j=1}^N$
- 12: For i = 1, ..., N set

$$v_j = c_j + \gamma \min_{a'} Q_{\theta_n}(s_{j+1}, a').$$

13: Let $L(\theta) := \sum_{j=1}^N |v_j - Q_{\theta}(s_j, a_j)|^2$ and update policy parameters

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} L(\theta).$$

14: end for

Classical policy gradient



Policy gradient (PG) methods

Recall we are minimizing

$$\Pi_{\mu} \ni \pi \mapsto V_{\tau}^{\pi}(\rho) \in \mathbb{R}$$
.

A "gradient" update would be

$$\pi_{n+1} = \pi_n - \eta \nabla_{\pi} V_{\tau}^{\pi}(\rho).$$

But even if S and A are of finite cardinality and

$$abla_{\pi}V_{ au}^{\pi}(
ho):=(
abla_{\pi(s,a)}V_{ au}^{\pi}(
ho))_{(s,a)\in S imes A}$$

with $(\nabla_{\pi(s,a)V_{\tau}^{\pi}(\rho)})_{(s,a)\in S\times A}\in \mathbb{R}^{N_S\times N_A}$ is a gradient in $\mathbb{R}^{N_S\times N_A}$ not in $\mathcal{P}(A|S)\equiv \Delta(A)^{N_S}$.

Policy gradient (PG) methods

$$\pi_{n+1} = \pi_n - \eta \nabla_\pi V_\tau^*(\rho).$$

Parametrize:

- ullet Direct: $rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(a|s) \propto e^{ heta(s,a)}$,
- Log-linear: $\frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s) \propto e^{\langle \theta, g(s,a) \rangle}$ with $g: S \times A \to \mathbb{R}^p$ basis.
- Neural-net: $\frac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(a|s) \propto e^{g_{ heta}(s,a)}$.

Then

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} V_{\tau}^{\pi_{\theta}}(\rho)$$

classical gradient descent: [Cauchy, 1847]¹ seems fine.

- **1** How to get $\nabla_{\theta} V_{\tau}^{\pi_{\theta}}(\rho)$ from data?
- ② Convergence: e.g. is $\theta \mapsto V_{\tau}^{\pi_{\theta}}(\rho)$ convex?

¹From [Lemaréchal, 2012] Cauchy and the gradient method, *Doc. Math. Extra*, 251–254.

Stochastic representation for solutions of certain linear equations

Let

$$d^{\pi}(ds'|s) = (1 - \gamma) \sum_{n=0}^{\infty} \gamma^{n} P_{\pi}^{n}(ds'|s) \quad \text{and} \quad d_{\rho}^{\pi}(ds) = \int_{S} d^{\pi}(ds|s') \rho(ds'). \tag{10}$$

We will refer to d^{π} as the occupancy kernel.

Lemma 5

Let $\pi \in \mathcal{P}(A|S)$ and $f, g \in B_b(S)$ such that for all $s \in S$,

$$f(s) = \gamma \int_{A} \int_{S} f(s') P(ds'|s, a) \pi(da|s) + g(s). \tag{11}$$

Then
$$f(s) = \frac{1}{1-\gamma} \int_S g(s') d^{\pi}(ds'|s)$$
 for all $s \in S$.

Proof of stochastic representation for solutions of certain linear equations

Proof. A kernel $k \in b\mathcal{M}(S|S)$ induces a linear operator $L_k \in \mathcal{L}(B_b(S))$ by

$$B_b(S) \ni h \mapsto L_k h = \int_S h(s') k(ds'|\cdot).$$

Since $||L_k h||_{B_b(S)} \le ||h||_{B_b(S)} ||k||_{b\mathcal{M}(S|S)}$ for all $h \in B_b(S)$, $||L_k||_{\mathcal{L}(B_b(S))} \le ||k||_{b\mathcal{M}(S|S)}$.

Consider the kernel $\gamma P_{\pi} \in b\mathcal{M}(S|S)$ defined by $(\gamma P_{\pi})(B) = \gamma \int_{B} \int_{A} P(ds'|s,a)\pi(da|s)$ for all $B \in \mathcal{B}(S)$. Then as $P_{\pi} \in \mathcal{P}(S|S)$ and $\|P_{\pi}\|_{b\mathcal{M}(S|S)} = 1$,

$$||L_{\gamma P_{\pi}}||_{\mathcal{L}(B_b(S))} \le ||\gamma P_{\pi}||_{b\mathcal{M}(S|S)} = \gamma ||P_{\pi}||_{b\mathcal{M}(S|S)} = \gamma < 1.$$

The linear equation (11) that f satisfies g is equivalent to

$$(id-L_{\gamma P_{\pi}})f=g$$
.

The operator id $-L_{\gamma P_{\pi}} \in \mathcal{L}(B_b(S))$ is invertible, and the inverse operator is given by the Neumann series

$$(\mathsf{id} - L_{\gamma P_\pi})^{-1} = \sum_{n=0}^\infty L_{\gamma P_\pi}^n$$
.

Thus, $f = \sum_{n=0}^{\infty} L_{\gamma P_{\pi}}^{n} g$. Observe that $L_{\gamma P_{\pi}}^{n} = L_{\gamma^{n} P_{\pi}^{n}}$ for all $n \in \mathbb{N}_{0}$, where P_{π}^{n} is the *n*-times product of the kernel P_{π} with $P_{\pi}^{0}(ds'|s) \coloneqq \delta_{s}(ds')$. Then by the definition (10) of $d^{\pi} \in \mathcal{P}(S|S)$,

$$f = \sum_{n=0}^{\infty} L_{\gamma P_{\pi}}^{n} g = \frac{1}{1-\gamma} \int_{S} g(s') d^{\pi}(ds'|\cdot)$$

which is the desired identity. \square

Performance difference

Lemma 6 (Performance difference¹)

For all $\rho \in \mathcal{P}(S)$ and $\pi, \pi' \in \Pi_{\mu}$,

$$egin{aligned} V^\pi_ au(
ho) - V^{\pi'}_ au(
ho) &= rac{1}{1-\gamma} \int_\mathcal{S} \int_A \Big(Q^{\pi'}_ au(s,a) + au \ln rac{\mathrm{d}\pi'}{\mathrm{d}\mu}(a|s) - V^{\pi'}_ au(s)\Big) (\pi-\pi')(da|s) d^\pi_
ho(ds) \ &+ rac{ au}{1-\gamma} \int_\mathcal{S} \mathsf{KL}(\pi(\cdot|s)|\pi'(\cdot|s)) d^\pi_
ho(ds) \,. \end{aligned}$$

¹Tabular case [Howard, 1960, Ch. 7, p. 87], re-discovered in RL context [Kakade and Langford, 2002], Polish spaces + entropy [Kerimkulov et al., 2025a]

Proof. By (7), for all $s \in S$,

$$\begin{split} &V_{\tau}^{\pi}(s) - V_{\tau}^{\pi'}(s) \\ &= \int_{A} \left(Q_{\tau}^{\pi}(a|s) + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a|s) \right) \pi(da|s) - \int_{A} \left(Q_{\tau}^{\pi'}(s,a) + \tau \ln \frac{\mathrm{d}\pi'}{\mathrm{d}\mu}(a|s) \right) \pi'(da|s) \\ &= \int_{A} \left(Q_{\tau}^{\pi'}(s,a) + \tau \ln \frac{\mathrm{d}\pi'}{\mathrm{d}\mu}(a|s) \right) (\pi - \pi') (da|s) \\ &+ \int_{A} \left(Q_{\tau}^{\pi}(s,a) + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a|s) - Q_{\tau}^{\pi'}(s,a) - \tau \ln \frac{\mathrm{d}\pi'}{\mathrm{d}\mu}(a|s) \right) \pi(da|s) \,. \end{split}$$

Hence for all $s \in S$ we have

$$egin{aligned} V^\pi_ au(s) - V^{\pi'}_ au(s) &= \int_A \left(Q^{\pi'}_ au(s, \mathsf{a}) + au \ln rac{\mathrm{d}\pi'}{\mathrm{d}\mu}(\mathsf{a}|s)
ight)(\pi - \pi')(d\mathsf{a}|s) \ &+ \gamma \int_A \int_S \left(V^\pi_ au(s') - V^{\pi'}_ au(s')
ight) P(ds'|s, \mathsf{a})\pi(d\mathsf{a}|s) + au \, \mathsf{KL}(\pi(\cdot|s)|\pi'(\cdot|s)) \,, \end{aligned}$$

where the last equality used def. of Q. fn (6) and KL identity (3). Hence, by Fubini's theorem and Lemma 5, for all $s \in S$,

$$\begin{split} &V_{\tau}^{\pi}(s) - V_{\tau}^{\pi'}(s) \\ &= \frac{1}{1-\gamma} \int_{S} \left[\int_{A} \left(Q_{\tau}^{\pi'}(s',a) + \tau \ln \frac{\mathrm{d}\pi'}{\mathrm{d}\mu}(a|s') \right) (\pi - \pi') (da|s') + \tau \operatorname{\mathsf{KL}}(\pi(\cdot|s')|\pi'(\cdot|s')) \right] d^{\pi}(ds'|s). \end{split}$$

Integrating both sides with respect to ho yields the desired identity. \Box

Towards PG for general state and action spaces

Proposition 4

Let
$$\pi, \pi' \in \Pi_{\mu}$$
 be such that $\pi(da|s) = \frac{\exp(f(s,a))\mu(da)}{\int_{A} \exp(f(s,a'))\mu(da')}$ for all $s \in S$. Then
$$\|Q_{\tau}^{\pi'} - Q_{\tau}^{\pi}\|_{B_{b}(S \times A)} \leq \frac{\gamma}{(1-\gamma)^{2}} \left(\|c\|_{B_{b}(S \times A)} + 2\tau \|f\|_{B_{b}(S \times A)}\right) \|\pi - \pi'\|_{b\mathcal{M}(A|S)} + \frac{\tau\gamma}{1-\gamma} \left\|\ln\frac{\mathrm{d}\pi'}{\mathrm{d}\pi}\right\|_{B_{b}(S \times A)}.$$

Proof. Start by getting the estimate for $\|V_{\tau}^{\pi'} - V_{\tau}^{\pi}\|_{B_b(S)}$ using Lemma 6 (performance difference).

PG for general state and action spaces

Proposition 5

Let $\tau \geq 0$ and $\rho \in \mathcal{P}(S)$. For all $\pi, \pi' \in \Pi_{\mu} \subset \mathcal{P}(A|S)$

$$\lim_{\varepsilon \searrow 0} \frac{V_{\tau}^{(1-\varepsilon)\pi+\varepsilon\pi'}(\rho) - V_{\tau}^{\pi}(\rho)}{\varepsilon} = \frac{1}{1-\gamma} \int_{\mathcal{C}} \int_{\mathcal{C}} \left(Q_{\tau}^{\pi}(s,a) + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(a|s) - V_{\tau}^{\pi}(s) \right) (\pi' - \pi)(da|s) d_{\rho}^{\pi}(ds) \,.$$
(12)

Proof. Let
$$\pi^{\varepsilon}=(1-\varepsilon)\pi+\varepsilon\pi'=\pi+\varepsilon(\pi'-\pi)$$
 and note that $\pi-\pi^{\varepsilon}=-\varepsilon(\pi'-\pi)=\varepsilon(\pi-\pi')$. Then
$$\frac{1}{\varepsilon}(V^{\pi}_{\tau}(\rho)-V^{\pi^{\varepsilon}}_{\tau}(\rho))=\frac{1}{\varepsilon}\frac{1}{1-\gamma}\int_{S}\int_{A}\left(Q^{\pi^{\varepsilon}}_{\tau}(s,a)+\tau\ln\frac{\mathrm{d}\pi^{\varepsilon}}{\mathrm{d}\mu}(a|s)\right)(\pi-\pi^{\varepsilon})(da|s)d^{\pi}_{\rho}(ds)\\ +\frac{1}{\varepsilon}\frac{\tau}{1-\gamma}\int_{S}\mathrm{KL}(\pi(\cdot|s)|\pi^{\varepsilon}(\cdot|s))d^{\pi}_{\rho}(ds)\\ =\frac{1}{1-\gamma}\int_{S}\int_{A}\left(Q^{\pi^{\varepsilon}}_{\tau}(s,a)+\tau\ln\frac{\mathrm{d}\pi^{\varepsilon}}{\mathrm{d}\mu}(a|s)\right)(\pi-\pi')(da|s)d^{\pi}_{\rho}(ds)\\ +\frac{1}{\varepsilon}\frac{\tau}{1-\gamma}\int_{S}\mathrm{KL}(\pi(\cdot|s)|\pi^{\varepsilon}(\cdot|s))d^{\pi}_{\rho}(ds)\,.$$

Proof of PG theorem for general state and action spaces

From the KL identity (3) we get

$$egin{aligned} rac{1}{arepsilon}(V^\pi_ au(
ho)-V^{\pi^arepsilon}_ au(
ho)) &= rac{1}{1-\gamma}\int_{\mathcal{S}}\int_{A}Q^{\pi^arepsilon}_ au(s,a)(\pi-\pi')(da|s)d^\pi_
ho(ds) \ &+rac{1}{arepsilon}rac{ au}{1-\gamma}\int_{\mathcal{S}}\Big(\operatorname{KL}(\pi(\cdot|s)|\mu(\cdot|s))-\operatorname{KL}(\pi^arepsilon(\cdot|s)|\mu(\cdot|s))\Big)d^\pi_
ho(ds)\,. \end{aligned}$$

Thus

$$egin{aligned} rac{1}{arepsilon}(V^{\pi^arepsilon}_{ au}(
ho) - V^\pi_{ au}(
ho)) &= rac{1}{1-\gamma}\int_{\mathcal{S}}\int_{A}Q^{\pi^arepsilon}_{ au}(s,a)(\pi'-\pi)(da|s)d^\pi_
ho(ds) \ &+ rac{ au}{1-\gamma}\int_{\mathcal{S}}rac{1}{arepsilon}\Big(\operatorname{KL}(\pi^arepsilon(\cdot|s)|\mu(\cdot|s)) - \operatorname{KL}(\pi(\cdot|s)|\mu(\cdot|s))\Big)d^\pi_
ho(ds)\,. \end{aligned}$$

The first integral on the right hand side converges to $\frac{1}{1-\gamma}\int_S\int_AQ^\pi_\tau(s,a)(\pi'-\pi)(da|s)d^\pi_\rho(ds)$ as $\varepsilon\to 0$ due to Proposition 4. Moreover, as $\pi,\pi'\in\Pi_\mu$, for all $s\in S$, by [Kerimkulov et al., 2025b, Lemma 3.8],

$$\lim_{\varepsilon\searrow 0}\frac{1}{\varepsilon}\Big(\operatorname{KL}(\pi^\varepsilon(\cdot|s)|\mu(\cdot|s))-\operatorname{KL}(\pi(\cdot|s)|\mu(\cdot|s))\Big)=\int_A\ln\frac{\mathrm{d}\pi}{\mathrm{d}\mu}(s|s)(\pi'-\pi)(ds|s)\,,$$

which along with Proposition 1 and the dominated yields the desired limit.

First variation and chain rule

For a fixed $\nu \in \mathcal{P}(S)$ define $\langle \cdot, \cdot \rangle_{\nu} : B_b(S \times A) \times b\mathcal{M}(A|S) \to \mathbb{R}$ by

$$\langle Z, m \rangle_{\nu} = \frac{1}{1-\gamma} \int_{S} \int_{A} Z(s, a) m(da|s) \nu(ds), \quad (Z, m) \in B_b(S \times A) \times b \mathcal{M}(A|S).$$

As a consequence of Proposition 5, given $\nu \in \mathcal{P}(S)$ satisfying $d_{\rho}^{\pi} \ll \nu$,

$$\lim_{\varepsilon \searrow 0} \frac{V_{\tau}^{(1-\varepsilon)\pi+\varepsilon\pi'}(\rho) - V_{\tau}^{\pi}(\rho)}{\varepsilon} = \left\langle \frac{\delta V_{\tau}^{\pi}(\rho)}{\delta\pi} \Big|_{\nu}, \pi' - \pi \right\rangle_{\nu},$$

with

$$\frac{\delta V_{\tau}^{\pi}(\rho)}{\delta \pi}\bigg|_{\nu}(s,a) = \left(Q_{\tau}^{\pi}(s,a) + \tau \ln \frac{\mathrm{d}\pi}{\mathrm{d}\mu}(s,a) - V_{\tau}^{\pi}(s)\right) \frac{\mathrm{d}d_{\rho}^{\pi}}{\mathrm{d}\nu}(s). \tag{13}$$

Let $(\mathbb{H}, (\cdot, \cdot)_{\mathbb{H}})$ be a Hilbert space.

Lemma 7 (Chain rule)

Let
$$\pi: \mathbb{H} \to \Pi_{\mu}$$
 be given. Then $\partial_{\theta_i} V_{\tau}^{\pi_{\theta}}(\rho) = \left\langle \frac{\delta V_{\tau}^{\pi}(\rho)}{\delta \pi}, \partial_{\theta_i} \pi_{\theta} \right\rangle_{d_{\pi}^{\pi}}$.

Proof. Similar to [Kerimkulov et al., 2025a, Proposition 3.8].

Policy gradient theorem

Theorem 8 (PG for parametrization)

Let $rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(\mathsf{a}|\mathsf{s}):=rac{\mathrm{e}^{g_{ heta}(\mathsf{s},\mathsf{a})}}{Z_{ heta}(\mathsf{s})}$, where $Z_{ heta}(\mathsf{s}):=\int_{\mathcal{A}}\mathrm{e}^{g_{ heta}(\mathsf{s},\mathsf{a}')}\mu(\mathsf{d}\mathsf{a}')$. Then

$$\left[
abla_{ heta}^{\pi_{ heta}}(
ho) = rac{1}{1-\gamma} \mathbb{E}_{a \sim \pi_{ heta}(\cdot|s)}^{s \sim \sigma_{
ho}^{\pi_{ heta}}} \left[rac{\delta V_{ au}^{\pi_{ heta}}}{\delta \pi}(s,a)
abla_{ heta} \ln rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \mu}(a|s)
ight].$$

Proof. From Lemma 7 (chain rule) we have:

$$abla_{ heta} V_{ au}^{\pi_{ heta}}(
ho) = rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{A} rac{\delta V_{ au}^{\pi_{ heta}}}{\delta \pi}(s,a)
abla_{ heta} rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \mu}(a|s) \mu(da) \, d_{
ho}^{\pi_{ heta}}(ds) \, .$$

Taking the gradient of the logarithm and re-arranging we see that

$$\nabla_{\theta} \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s) = \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s)\nabla_{\theta} \ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s). \tag{14}$$

Hence

$$abla_{ heta} V_{ au}^{\pi_{ heta}}(
ho) = rac{1}{1-\gamma} \int_{\mathbb{R}} \int_{\mathbb{R}} rac{\delta V_{ au}^{\pi_{ heta}}}{\delta \pi}(s,a)
abla_{ heta} \ln rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \mu}(a|s) \pi_{ heta}(da|s) \, d_{
ho}^{\pi_{ heta}}(ds) \, .$$

We just need to rewrite this in terms of expectation to get the conclusion.

Some remarks of PG

Remark on baseline. We can take any $b \in B_b(S)$. Then

$$egin{aligned} &\int_A b(s)
abla_ heta \ln rac{\mathrm{d}\pi_ heta}{\mathrm{d}\mu}(a|s)\pi_ heta(da|s) = b(s) \int_A
abla_ heta \ln rac{\mathrm{d}\pi_ heta}{\mathrm{d}\mu}(a|s)\pi_ heta(da|s) = b(s) \int_A
abla_ heta rac{\mathrm{d}\pi_ heta}{\mathrm{d}\mu}(a|s)\mu(da) \\ &= b(s)
abla_ heta \int_A rac{\mathrm{d}\pi_ heta}{\mathrm{d}\mu}(a|s)\mu(da) = b(s)
abla_ heta 1 = 0 \ . \end{aligned}$$

Hence

$$abla_{ heta} V_{ au}^{\pi_{ heta}}(
ho) = rac{1}{1-\gamma} \mathbb{E}_{m{a} \sim \pi_{ heta}(\cdot|m{s})}^{s \sim d_{
ho}^{\pi_{ heta}}} \Big[\Big(rac{\delta V_{ au}^{\pi_{ heta}}}{\delta \pi}(m{s},m{a}) + b(m{s}) \Big)
abla_{ heta} \ln rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \mu}(m{a}|m{s}) \Big] \,.$$

Some remarks on PG

Remark on estimating the advantage function. First variation:

$$rac{\delta V_{ au}^{\pi \, heta}(s, a)}{\delta \pi}(s, a) = \underbrace{Q_{ au}^{\pi}(s, a) - V_{ au}^{\pi}(s)}_{=:A_{ au}^{\pi}(s, a) \text{ "advantage"}} + au \ln rac{\mathrm{d} \pi}{\mathrm{d} \mu}(s, a) \, .$$

Advantage $A_{\tau}^{\pi}(s, a)$ can be estimated from data: $(s_t, a_t, c_t, s_{t+1}, a_{t+1}, \ldots)$.

$$\hat{A}^\pi_ au := c_t + \gamma \hat{V}^\pi_ au(s_{t+1}) - \hat{V}^\pi_ au(s_t),$$

where $\hat{V}_{\tau} \approx V_{\tau}^{\pi}$. N.B.

$$\mathbb{E}_{s_{t+1} \sim P(\cdot|s_t, a_t)} \hat{A}^\pi_\tau = c(s_t, a_t) + \gamma \int_S \hat{V}_\tau(s') P(ds'|s_t, a_t) - \hat{V}_\tau(s_t)$$

would be equal to $A_{\tau}^{\pi}(s_t,a_t)$ if $\hat{V}_{\tau}=V_{\tau}^{\pi}$ in which case it would be unbiased. Alternative

$$\hat{A}^\pi_ au := \sum_{l=0}^\infty \gamma^{t+l} c_{t+l} - \hat{V}(\mathsf{s}_t)\,.$$

Generalised advantage estimation (GAE) formula [Schulman et al., 2015] allows efficient variance vs bias tradeoffs.

Corollary 9 (to Policy Gradient Theorem)

Let $rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(\mathsf{a}|\mathsf{s}):=rac{\mathrm{e}^{\mathsf{g}_{ heta}(\mathsf{s},\mathsf{a})}}{\mathsf{Z}_{ heta}(\mathsf{s})}$, $\mathsf{Z}_{ heta}(\mathsf{s}):=\int_{\mathsf{A}}\mathsf{e}^{\mathsf{g}_{ heta}(\mathsf{s},\mathsf{a}')}\mu(\mathsf{d}\mathsf{a}')$. Then

$$abla_{ heta} V^{\pi_{ heta}}_{ au}(
ho) = rac{1}{1-\gamma} \mathbb{E}^{s\sim d^{\pi_{ heta}}_{
ho}}_{a\sim \pi_{ heta}(\cdot|s)} iggl[rac{\delta V^{\pi_{ heta}}_{ au}}{\delta \pi}(s,a) iggl(
abla_{ heta} g_{ heta}(s,a) - \int_{A} (
abla_{ heta} g_{ heta})(s,a') \pi_{ heta}(da'|s) iggr) iggr] \,.$$

Proof. Note that

$$\ln rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(a|s) = g_{ heta}(s,a) - \ln Z_{ heta}(s)$$

and so

$$abla_{ heta} \ln rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu}(s,a) =
abla_{ heta} g_{ heta}(s,a) -
abla_{ heta} Z_{ heta}(s) rac{1}{Z_{ heta}(s)} =
abla_{ heta} g_{ heta}(s,a) - \int_{A} (
abla_{ heta} g_{ heta})(s,a') rac{e^{g_{ heta}(s,a')}}{Z_{ heta}(s)} \mu(da') \,.$$

Hence we have an expression for gradient of the log-density:

$$\nabla_{\theta} \ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(s,a) = \nabla_{\theta} g_{\theta}(s,a) - \int_{A} (\nabla_{\theta} g_{\theta})(s,a') \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a'|s)\mu(da') \tag{15}$$

which concludes the calculation. \square

Some remarks on PG

If the state and action spaces are finite and we take the direct (tabular) parametrizations so that $g_{\theta}(s,a) := \theta(s,a)$ then

$$egin{aligned} \partial_{ heta_{\hat{s},\hat{s}}} g_{ heta}(s,a) &- \sum_{a'} \partial_{ heta_{\hat{s},\hat{s}}} g_{ heta}(s,a') \pi_{ heta}(a'|s) = \delta_{\hat{s},\hat{s}}(s,a) - \sum_{a'} \delta_{\hat{s},\hat{s}}(s,a') \pi(a'|s) \ &= \delta_{\hat{s},\hat{s}}(s,a) - \delta_{\hat{s}}(s) \pi(\hat{a}|s) = \delta_{\hat{s}}(s) (\delta_{\hat{s}}(a) - \delta_{\hat{s}}(s) \pi(\hat{a}|s)) \ . \end{aligned}$$

Hence

$$egin{aligned} \partial_{ heta_{\hat{s},\hat{s}}} V^{\pi_{ heta}}_{ au}(
ho) &= rac{1}{1-\gamma} \sum_{s,a} rac{\delta V^{\pi_{ heta}}_{ au}}{\delta \pi}(s,a) \delta_{\hat{s}}(s) \delta_{\hat{s}}(a) \pi_{ heta}(a|s) d^{\pi_{ heta}}_{
ho}(s) \ &- rac{1}{1-\gamma} \sum_{s,a} rac{\delta V^{\pi_{ heta}}_{ au}}{\delta \pi}(s,a) \delta_{\hat{s}}(s) \pi(\hat{a}|s) \pi_{ heta}(a|s) d^{\pi_{ heta}}_{
ho}(s) \,. \end{aligned}$$

But

$$\sum_{s,a} \frac{\delta V_{\tau}^{\pi_{\theta}}}{\delta \pi}(s,a) \delta_{\hat{s}}(s) \pi(\hat{a}|s) \pi_{\theta}(a|s) d_{\rho}^{\pi_{\theta}}(s) = \sum_{s} \delta_{\hat{s}}(s) \pi(\hat{a}|s) \sum_{a} \frac{\delta V_{\tau}^{\pi_{\theta}}}{\delta \pi}(s,a) \pi_{\theta}(a|s) d_{\rho}^{\pi_{\theta}}(s) = 0$$

and so

$$\partial_{ heta_{\hat{s},\hat{a}}} V^{\pi_{ heta}}_{ au}(
ho) = rac{1}{1-\gamma} rac{\delta V^{\pi_{ heta}}_{ au}}{\delta \pi} (\hat{s},\hat{a}) \pi_{ heta} (\hat{a}|\hat{s}) d^{\pi_{ heta}}_{
ho} (\hat{s}) \,.$$

This is (for the $\tau=0$ case) exactly Lemma C.1 in [Agarwal et al., 2019].

Some remarks on PG

If
$$g_{ heta}(s,a)=(heta,\phi(s,a))_{\mathbb{H}}$$
 then $\partial_{ heta_i}g_{ heta}(s,a)= heta_i(s,a)$ and so

$$\nabla_{\theta} V_{\tau}^{\pi_{\theta}}(\rho) = \frac{1}{1 - \gamma} \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}(\cdot|\mathbf{s})}^{\mathbf{s} \sim d_{\rho}^{\pi_{\theta}}} \left[\mathbf{s}, \mathbf{a} \right) \left(\phi(\mathbf{s}, \mathbf{a}) - \int_{\mathcal{A}} \phi(\mathbf{s}, \mathbf{a}') \pi_{\theta}(d\mathbf{a}'|\mathbf{s}) \right) \right]. \tag{16}$$

Summary of PG so far:

- We have expression for the gradient.
- It can be estimated from data.
- For some simple parametrizations it's nice and simple.

Next: what about convergence?

Lack of convexity in softmax parametrization

Consider *minimizing*, over $\pi \in \mathcal{P}(A)$ the objective

$$v^{\pi}:=\int_{A}c(a)\pi(da)$$
 .

We trivially have, for $\pi, \pi' \in \mathcal{P}(A)$ that

$$v^{(1-arepsilon)\pi+arepsilon\pi'} \leq (1-arepsilon)v^\pi + arepsilon v^{\pi'}$$

so $\mathcal{P}(A) \ni \pi \mapsto v^{\pi} \in \mathbb{R}$ is convex.

Consider *minimizing*, over $\theta \in \mathbb{R}^{|A|}$ the objective

$$v^{\pi_{ heta}}:=\int_{A}c(a)\pi_{ heta}(da),$$

with $\pi_{ heta}(a) = rac{e^{ heta(a)}}{\sum_{a'} e^{ heta(a')}}$.

The map $\mathbb{R}^p \ni \theta \mapsto v^{\pi_{\theta}} \in \mathbb{R}$ is *not* convex [Mei et al., 2020, Propn. 1].

Performance difference vs convexity

Definition (Convexity)

If for some $\tau \geq 0$ we have all $m, m' \in \mathcal{P}(A)$ that

$$F(m) - F(m') \ge \left\langle rac{\delta F(m')}{\delta m}, m - m'
ight
angle + au \, \mathsf{KL}(m|m') \, ,$$

then F is convex $(\tau = 0)$ or strongly convex $(\tau > 0)$.

Equiv.:
$$\frac{\delta F}{\delta m}(m,\cdot)$$
 exists and $F((1-\varepsilon)m+\varepsilon m')\leq (1-\varepsilon)F(m)+\varepsilon F(m')$ for all $m,m'\in\mathcal{P}(A),\ \varepsilon\in[0,1]$.

Performance difference

$$(V^\pi_{ au}-V^{\pi'}_{ au})(
ho)=\left\langlerac{\delta V^{\pi'}_{ au}}{\delta\pi},\pi-\pi'
ight
angle_{
ho,\pi}+rac{ au}{1-\gamma}\int_{\mathcal{S}}\mathsf{KL}(\pi|\pi')(s)d^\pi_
ho(ds)\,,$$

where

$$\langle h, \hat{\pi}
angle_{
ho,\pi} := rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{A} h(s,a) \hat{\pi}(da|s) \, d^{\pi}_{
ho}(ds)$$

The map $\Pi_{\mu} \ni \pi \mapsto V_{\tau}^{\pi}(\rho) \in \mathbb{R} \cup \{+\infty\}$ is **not** convex, e.g. [Giegrich et al., 2024, Proposition 2.4] **even if** underlying dynamics is linear and costs convex.

Polyak-Lojasiewicz condition: Gradient dominance

Continuous time gradient flow

$$\frac{d}{ds}\theta_s = -\nabla f(\theta_s) \implies \frac{d}{ds} \big[f(\theta_s) - f(\theta^*) \big] = -|\nabla f(\theta_s)|^2$$
.



Non-uniform Polyak–Łojasiewicz: there is $\mu:\mathbb{R}^p o (0,\infty)$ s.t. for all $\theta\in\mathbb{R}^p$

$$0 \le f(\theta) - f(\theta^*) \le \mu(\theta) |\nabla f(\theta)|^2$$
.

Hence

$$\frac{d}{ds}\big[f(\theta_s) - f(\theta^*)\big] = -|\nabla f(\theta_s)|^2 \le -\mu^{-1}(\theta_s)\big[f(\theta_s) - f(\theta^*)\big]$$

Grönwall:

$$0 \leq f(\theta_s) - f(\theta^*) \leq \left[f(\theta_0) - f(\theta^*) \right] \exp\left(- \int_0^s \mu^{-1}(\theta_r) \, dr \right).$$

Q: Is
$$\inf_r \mu^{-1}(\theta_r) \ge \alpha > 0$$
?

Convergence of classical PG in MDPs

- Discrete time LQR: Polyak–Łojasiewicz (PL) / gradient dominance established and so PG has linear convergence [Fazel et al., 2018, Bu et al., 2019, Hu et al., 2023].
- In general discrete state-action setting best PL result is non-uniform [Mei et al., 2020] but shown lower bounded along PG and hence convergence.

Mirror descent



Mirror descent

Static optimization mirror descent:

- Goes back to at least [Nemirovski, 1979].
- Modern proximal point form [Beck and Teboulle, 2003].
- For general probability measures [Aubin-Frankowski et al., 2022].

Mirror descent for MDPs

Discrete space MDPs and constants **dependent** on |S| and |A|:

- [Cen et al., 2022], entropy regularised, show linear convergence for disc. time. mirror descent
- [Cayci et al., 2021] same setting i.e natural policy gradient, log-linear policies i.e. mirror desc with func. approx.
- [Xiao, 2022] and [Khodadadian et al., 2022] achieved *linear convergence for unregularised MDPs* with inexact policy evaluation by employing geometrically increasing step sizes in NPG.

Discrete space MDPs and constants **independent of** of |S| and |A|:

• [Lan, 2023] linear convergence of policy mirror descent with arbitrary convex regularisers and [Zhan et al., 2023] convergence rates independent of action space dimension.

MDPs with **general** S and A:

• Discrete step mirror descent and Fisher-Rao flow: Exponential convergence for entropy regularized MDPs in Polish state & action spaces [Kerimkulov et al., 2025a].

Deriving mirror descent

Aim: find

$$\pi^*(\cdot|s) = \arg\min_{\pi} V^{\pi}_{\tau}(s)$$
.

Let's say we have π_{old} . Fix $\rho \in \mathcal{P}(S)$ and write $V_{\tau}^{\pi} = V_{\tau}^{\pi}(\rho)$. By perf. diff., Lemma 6,

$$V^\pi_ au = V^{\pi_{\mathsf{old}}}_ au + igl\langle rac{\delta V^{\pi_{\mathsf{old}}}_ au}{\delta \pi}, \pi - \pi_{\mathsf{old}} igr
angle_{
ho, oldsymbol{\pi}} + rac{ au}{1-\gamma} \int_{\mathcal{S}} \mathsf{KL}(\pi | \pi_{\mathsf{old}})(s) d^{oldsymbol{\pi}}_
ho(\mathit{ds}) \,.$$

Linearize and penalize with $\lambda \geq \tau$ to not move too far

$$L^{\pi} := V^{\pi_{\mathsf{old}}}_{ au} + \left\langle rac{\delta V^{\pi_{\mathsf{old}}}_{ au}}{\delta \pi}, \pi - \pi_{\mathsf{old}}
ight
angle_{
ho, \pi_{\mathsf{old}}} + rac{\lambda}{1 - \gamma} \int_{S} \mathsf{KL}(\pi | \pi_{\mathsf{old}})(s) d^{\pi_{\mathsf{old}}}_{
ho}(ds) \,.$$

Mirror descent optimizes $\pi \mapsto L^{\pi}(x)$ giving

$$\pi_{\mathsf{new}}(\mathit{da}|s) = \arg\min_{\pi} \left(V_{\tau}^{\pi_{\mathsf{old}}} + \int_{A} \frac{\delta V_{\tau}^{\pi_{\mathsf{old}}}}{\delta \pi}(s, a) (\pi - \pi_{\mathsf{old}}) (\mathit{da}|s) + \lambda \, \mathsf{KL}(\pi | \pi_{\mathsf{old}})(s) \right).$$

Motivation for studying mirror descent

Policy gradient: introduce parametrized densities $\pi_{\theta}(\mathit{da},s) \propto \mathrm{e}^{\mathrm{g}_{\theta}(a,s)}\mu(\mathit{da})$. Step $\eta > 0$:

$$egin{align*} heta_{\mathsf{new}} &= heta_{\mathsf{old}} + \eta
abla_{ au} V_{ au}^{\pi_{m{ heta}}_{\mathsf{old}}} \ , \ &
abla_{ au} V_{ au}^{\pi_{m{ heta}}_{\mathsf{old}}}(
ho) = rac{1}{1-\gamma} \mathbb{E}_{a \sim \pi_{m{ heta}_{\mathsf{old}}}(\cdot|s)}^{s \sim d_{
ho}^{\pi_{m{ heta}}_{\mathsf{old}}}} \left[rac{\delta V_{ au}^{\pi_{m{ heta}}_{\mathsf{old}}}}{\delta \pi}(s, a)
abla_{m{ heta}} \ln rac{\mathrm{d} \pi_{m{ heta}_{\mathsf{old}}}}{\mathrm{d} \mu}(a|s)
ight]. \end{split}$$

Problem: Even if θ_{new} and θ_{old} are close $\pi_{\theta_{\text{old}}}$ and $\pi_{\theta_{\text{new}}}$ may be very different!

Instead, re-write the mirror descent objective:

$$\begin{split} L_{\mathsf{MD}}(\theta) &= \left\langle \frac{\delta V_{\tau}^{\pi_{\theta_{\mathsf{old}}}}}{\delta \pi}, \pi_{\theta} \right\rangle_{\rho, \pi_{\theta_{\mathsf{old}}}} + \lambda \int_{S} \mathsf{KL}(\pi_{\theta} | \pi_{\theta_{\mathsf{old}}})(s) d_{\rho}^{\pi_{\theta_{\mathsf{old}}}}(ds) \\ &= \mathbb{E}_{\mathsf{a} \sim \pi_{\theta}(\cdot | s)}^{s \sim \theta_{\mathsf{old}}^{\theta_{\mathsf{old}}}} \left[\frac{\delta V_{\tau}^{\pi_{\theta_{\mathsf{old}}}}}{\delta \pi}(s, a) + \lambda \, \mathsf{KL}(\pi_{\theta} | \pi_{\theta_{\mathsf{old}}})(s) \right] \\ &= \mathbb{E}_{\mathsf{a} \sim \pi_{\theta_{\mathsf{old}}}^{\theta_{\mathsf{old}}}(\cdot | s)}^{s \sim \theta_{\mathsf{old}}^{\theta_{\mathsf{old}}}} \left(s, a \right) \frac{\mathrm{d} \pi_{\theta}}{\mathrm{d} \pi_{\theta_{\mathsf{old}}}}(a | s) + \lambda \, \mathsf{KL}(\pi_{\theta} | \pi_{\theta_{\mathsf{old}}})(s) \right]. \end{split}$$

Step $\eta > 0$:

$$heta_{\mathsf{new}} = heta_{\mathsf{old}} + \eta
abla_{ heta} L_{\mathsf{MD}}(heta_{\mathsf{old}})$$
 .

Mirror descent policy improvement (with exact update)

Mirror descent update

$$\pi^{n+1}(\cdot|s) = \underset{m \in \mathcal{P}(A)}{\operatorname{argmin}} \int_{A} \frac{\delta V_{\tau}^{\pi_{n}}}{\delta \pi}(s, a) (m(da) - \pi^{n}(da|s)) + \lambda \operatorname{KL}(m|\pi^{n}(\cdot|s)). \tag{17}$$

From the performance difference lemma, see Lemma (6), we see that

$$(V_{\tau}^{n+1} - V_{\tau}^{n})(\rho) = \frac{1}{1-\gamma} \int_{S} \left(\int_{A} \frac{\delta V_{\tau}^{n}}{\delta \pi}(s, a) (\pi^{n+1} - \pi^{n}) (da|s) + \tau \operatorname{KL}(\pi^{n+1}|\pi^{n})(s) \right) d_{\rho}^{\pi^{n+1}}(ds)$$

$$\leq \frac{1}{1-\gamma} \int_{S} \left(\int_{A} \frac{\delta V_{\tau}^{n}}{\delta \pi}(s, a) (\pi^{n+1} - \pi^{n}) (da|s) + \lambda \operatorname{KL}(\pi^{n+1}|\pi^{n})(s) \right) d_{\rho}^{\pi^{n+1}}(ds) .$$
(18)

From the mirror descent update (17) we have, for all $\pi \in \Pi_{\mu}$ and $s \in S$ that

$$\textstyle \int_A \frac{\delta V_\tau^n}{\delta \pi}(\mathsf{s},\mathsf{a})(\pi-\pi^n)(\mathsf{d}\mathsf{a}|\mathsf{s}) + \lambda \, \mathsf{KL}(\pi|\pi^n)(\mathsf{s}) \geq \int_A \frac{\delta V_\tau^n}{\delta \pi}(\mathsf{s},\mathsf{a})(\pi^{n+1}-\pi^n)(\mathsf{d}\mathsf{a}|\mathsf{s}) + \lambda \, \mathsf{KL}(\pi^{n+1}|\pi^n)(\mathsf{s}) \,.$$

This with $\pi = \pi^n$ allows us to conclude that for all $s \in S$ we have

$$\int_{A} \frac{\delta V_{n}^{r}}{\delta \pi}(s,a)(\pi^{n+1} - \pi^{n})(da|s) + \lambda \operatorname{KL}(\pi^{n+1}|\pi^{n})(s) \leq 0.$$
(19)

From (18) we have

$$(V_{\tau}^{n+1}-V_{\tau}^n)(\rho)\leq 0.$$

Mirror descent with approximation

Recall that
$$rac{\delta V_{ au}^{\pi_n}}{\delta \pi} = A_{ au}^{\pi_n} + au \ln rac{\mathrm{d} \pi^n}{\mathrm{d} \mu} = Q_{ au}^{\pi_n} - V_{ au}^{\pi_n} + au \ln rac{\mathrm{d} \pi^n}{\mathrm{d} \mu}.$$

Updates can only be made with an approximation $\hat{A}_n(s,a) = A_{\tau}^{\pi_n}(s,a) + \mathcal{E}_n(s,a)$.

Consider the scheme

$$\pi^{n+1}(da|s) = \underset{m \in \mathcal{P}(A)}{\operatorname{argmin}} \int_{A} \left(\hat{A}_{n}(s,a) + \tau \ln \frac{\mathrm{d}\pi^{n}}{\mathrm{d}\mu}(a|s) \right) (m(da) - \pi^{n}(da|s)) + \lambda \operatorname{KL}(m|\pi^{n}(\cdot|s)). \tag{20}$$

Towards L-smoothness

This is from [Lan, 2023].

Lemma 10

Let $F: S \to \mathbb{R}$ be such that $F \leq 0$. Then for any π and any $s \in S$

$$\frac{1}{1-\gamma}\int_{S}F(s')\,d_{s}^{\pi}(ds')\leq F(s). \tag{21}$$

Proof. From (10) and the fact that $P_{\pi}^{0}(ds'|s) = \delta_{s}(ds')$ we have for all $s \in S$ that

$$\frac{1}{1-\gamma} \int_{S} F(s') d_{s}^{\pi}(ds') = \int_{S} F(s') P_{\pi}^{0}(ds'|s) + \sum_{k=1}^{\infty} \int_{S} \gamma^{k} F(s') P_{\pi}^{k}(ds'|s)
\leq \int_{S} F(s') \delta_{s}(ds') = F(s).$$
(22)

This concludes the proof. \Box

Lemma 11 (L-smoothness for exact update)

Let $\pi, \pi' \in \Pi_{\mu}$ satisfy $\int_{A} \frac{\delta V_{\pi}^{\pi'}}{\delta \pi}(s, a)(\pi - \pi')(da|s) + \tau \operatorname{KL}(\pi|\pi')(s) \leq 0$ for all $s \in S$. Then for all $s \in S$,

$$(V^\pi_ au-V^{\pi'}_ au)(s) \leq \int_{\mathbb{A}} rac{\delta \, V^{\pi'}_ au}{\delta \pi}(s,\mathsf{a})(\pi-\pi')(\mathsf{d}\mathsf{a}|s) + au\,\mathsf{KL}(\pi|\pi')(s)\,.$$

In particular with $\pi' = \pi_{old}$ and $\pi = \pi_{new}$ given by the exact update (17) satisfy this.

Proof. From perf. diff. lemma and Lan's trick:

$$egin{aligned} (V^\pi_ au - V^{\pi'}_ au')(s) &\leq rac{1}{1-\gamma} \int_S ig(\int_A rac{\delta V^{\pi'}_ au'}{\delta \pi}(s',a)(\pi-\pi')(da|s') + au \, \mathsf{KL}(\pi|\pi')(s') ig) d^\pi_s(ds') \ &\leq \int_A rac{\delta V^{\pi'}_ au'}{\delta \pi}(s,a)(\pi-\pi')(da|s) + au \, \mathsf{KL}(\pi|\pi')(s) \,. \end{aligned}$$

Convergence of mirror descent with approximate advantage



Ingredients for convergence of mirror descent

- Convexity (strong for "linear" rate)
- L-smoothness
- Three point lemma

Lemma 12 (Three point lemma / Bregman proximal inequality)

Let $G:M_{\mu} \to \mathbb{R}$ be convex. For all $m' \in M_{\mu}$ let

$$m^* = \operatorname*{argmin}_{m \in M_u} \left\{ G(m) + \mathsf{KL}(m|m') \right\} . \tag{23}$$

Then, for all $m \in M_{\mu}$, we have

$$G(m) + KL(m|m') \ge G(m^*) + KL(m|m^*) + KL(m^*|m').$$
 (24)

The proof of Lemma 12 can be found e.g., in [Aubin-Frankowski et al., 2022] noting that the flat derivative of KL is well defined on M_{μ} , see e.g. [Kerimkulov et al., 2025b, Lemma 3.8].

Let π^n be generated by inductive application of the approximate mirror descent step (20). Let $V^n_{\tau} := V^{\pi^n}_{\tau}$ for $n \in \mathbb{N}$. We begin with an application of Bregman proximal inequality, see Lemma 12. Fix $s \in S$ and $\pi^n \in \Pi_{\mu}$ and define $G: M_{\mu} \to \mathbb{R}$ by

$$G(m) = rac{1}{\lambda} \int_{A} \Big(\hat{A}_n(s,a) + au \ln rac{\mathrm{d} \pi^n}{\mathrm{d} \mu}(a|s) \Big) (m(da) - \pi^n(da|s)) \,.$$

It is linear and thus clearly convex and hence due to the mirror descent update (20) is equivalent to (23) and so we have, for all $\pi \in \Pi_{\mu}$, $s \in S$ and $n \in \mathbb{N}$ that

$$egin{split} & rac{1}{\lambda} \int_{A} \Big(\hat{A}_n(s,a) + au \ln rac{\mathrm{d}\pi^n}{\mathrm{d}\mu}(a|s)\Big) (\pi-\pi^n)(da|s) + \mathsf{KL}(\pi|\pi^n)(s) \ & \geq rac{1}{\lambda} \int_{A} \Big(\hat{A}_n(s,a) + au \ln rac{\mathrm{d}\pi^n}{\mathrm{d}\mu}(a|s)\Big) (\pi^{n+1}-\pi^n)(da|s) + \mathsf{KL}(\pi|\pi^{n+1})(s) + \mathsf{KL}(\pi^{n+1}|\pi^n)(s) \,. \end{split}$$

Re-arranging this leads to

$$KL(\pi|\pi^{n+1})(s) - KL(\pi|\pi^{n})(s)
\leq \frac{1}{\lambda} \int_{A} \left(\hat{A}_{n}(s,a) + \tau \ln \frac{d\pi^{n}}{d\mu}(a|s) \right) (\pi - \pi^{n})(da|s)
- \frac{1}{\lambda} \int_{A} \left(\hat{A}_{n}(s,a) + \tau \ln \frac{d\pi^{n}}{d\mu}(a|s) \right) (\pi^{n+1} - \pi^{n})(da|s) - KL(\pi^{n+1}|\pi^{n})(s) .$$
(25)

From the performace difference, Lemma 6, we have

$$\begin{split} &(V_{\tau}^{n+1}-V_{\tau}^{n})(s)\\ &=\frac{1}{1-\gamma}\int_{\mathcal{S}}\big(\int_{A}\big(\hat{A}_{n}-\mathcal{E}_{n}+\tau\ln\frac{\mathrm{d}\pi^{n}}{\mathrm{d}\mu}\big)(s,a)(\pi^{n+1}-\pi^{n})(da|s)+\tau\operatorname{\mathsf{KL}}(\pi^{n+1}|\pi^{n})(s)\big)d_{\rho}^{\pi^{n+1}}(ds)\,. \end{split}$$

Note that (20), together with $\lambda \geq \tau$ guarantees that

$$0 \geq \int_{\mathcal{A}} ig(\hat{A}_n(s,a) + au \ln rac{\mathrm{d}\pi^n}{\mathrm{d}\mu}(a|s)ig) ig(\pi^{n+1} - \pi^n)(da|s) + au \operatorname{\mathsf{KL}}(\pi^{n+1}|\pi^n)(s) =: \mathcal{F}(s)$$

for all $s \in S$. Thus we may apply Lemma 10 and get

$$(V_{\tau}^{n+1}-V_{\tau}^{n})(s) \leq F(s)-\frac{1}{1-\gamma}\int_{S}\int_{A}\mathcal{E}_{n}(s,a)(\pi^{n+1}-\pi^{n})(da|s)d_{\rho}^{\pi^{n+1}}(ds)$$
.

Assume that $\|\mathcal{E}\|_{B_h(S\times A)}=\delta_n<\infty$. Then we have the following approximate L-smoothness:

$$(V_{\tau}^{n+1}-V_{\tau}^n)(s)\leq F(s)+rac{2\delta_n}{1-\gamma}\,,\;\;s\in S.$$

Applying this in (25) and taking we thus have, for all $s \in S$, that

$$\mathsf{KL}(\pi_{\tau}^{*}|\pi^{n+1})(s) - \mathsf{KL}(\pi_{\tau}^{*}|\pi^{n})(s) \leq \frac{1}{\lambda} \int_{A} \left(\hat{A}_{n}(s,a) + \tau \ln \frac{\mathrm{d}\pi^{n}}{\mathrm{d}\mu}(a|s)\right) (\pi_{\tau}^{*} - \pi^{n})(da|s) - \frac{1}{\lambda} (V_{\tau}^{n+1} - V_{\tau}^{n})(s) + \frac{2\delta_{n}}{(1-\gamma)\lambda}.$$
(26)

Summing up over $n=0,1,\ldots,N-1$ we see (spotting the telescoping sums) that for all $s\in S$,

$$\mathsf{KL}(\pi_ au^*|\pi^N)(s) - \mathsf{KL}(\pi_ au^*|\pi^0)(s) \leq \sum_{n=0}^{N-1} rac{1}{\lambda} \int_{\mathcal{A}} \Big(\hat{A}_n(s,a) + au \ln rac{\mathrm{d}\pi^n}{\mathrm{d}\mu}(a|s)\Big) (\pi_ au^* - \pi^n)(da|s) \ - rac{1}{\lambda} (V_ au^N - V_ au^0)(s) + rac{2}{(1-\gamma)\lambda} \sum_{r=0}^{N-1} \delta_n \,.$$

We wish to apply performance difference in due course and so we observe that the above is equivalent to

$$\mathsf{KL}(\pi_{\tau}^{*}|\pi^{N})(s) - \mathsf{KL}(\pi_{\tau}^{*}|\pi^{0})(s) \leq \sum_{n=0}^{N-1} \frac{1}{\lambda} \int_{A} \left(A_{\tau}^{\pi^{n}}(s,a) + \tau \ln \frac{\mathrm{d}\pi^{n}}{\mathrm{d}\mu}(a|s) \right) (\pi_{\tau}^{*} - \pi^{n})(da|s) + \sum_{n=0}^{N-1} \frac{1}{\lambda} \int_{A} \mathcal{E}_{n}(s,a) (\pi_{\tau}^{*} - \pi^{n})(da|s) - \frac{1}{\lambda} (V_{\tau}^{N} - V_{\tau}^{0})(s) + \frac{2}{(1-\gamma)\lambda} \sum_{n=0}^{N-1} \delta_{n}.$$
(27)

Notice that $V^N_{\tau}(s) \geq V^*_{\tau}(s)$ and so $(V^N_{\tau} - V^0_{\tau})(s) \geq (V^*_{\tau} - V^0_{\tau})(s)$ for all $N \in \mathbb{N}$. Let

$$y^n:=\int_{\mathcal{S}}\mathsf{KL}(\pi^*_ au|\pi^n)(s)d^{\pi^*_ au}_
ho(ds)$$
 and $lpha:=-\int_{\mathcal{S}}(V^*_ au-V^0)(s)d^{\pi^*_ au}_
ho(ds)$

so that, after integrating (27) over $d_{\rho}^{\pi_{\tau}^*}$ and using $\|\mathcal{E}\|_{B_{\delta}(S\times A)} = \delta_n < \infty$ we have

$$y^N-y^0 \leq \sum_{r=0}^{N-1} \frac{1}{\lambda} \int_S \int_A \frac{\delta V_\tau^n}{\delta \pi}(s,a) (\pi_\tau^*-\pi^n) (\mathsf{d} \mathsf{a} | s) \mathsf{d}_\rho^{\pi_\tau^*}(\mathsf{d} s) + \frac{2}{\lambda} \sum_{r=0}^{N-1} \delta_n + \frac{\alpha}{\lambda} + \frac{2}{(1-\gamma)\lambda} \sum_{r=0}^{N-1} \delta_n \,.$$

Using the performance difference lemma, see Lemma 6, and upper bounding the approximation error terms we get

$$y^{N}-y^{0} \leq \sum_{n=0}^{N-1} \left[\frac{1-\gamma}{\lambda} (V^{\pi^*_{\tau}}-V^{\pi^n})(\rho) - \frac{\tau}{\lambda} \int_{\mathcal{S}} \mathsf{KL}(\pi^*_{\tau}|\pi^n)(s) d^{\pi^*_{\tau}}_{\rho}(ds)\right] + \frac{\alpha}{\lambda} + \frac{4}{(1-\gamma)\lambda} \sum_{n=0}^{N-1} \delta_n.$$

Since since $\mathsf{KL}(\cdot|\cdot) \geq 0$ we get that

$$y^N - y^0 \leq N \frac{1-\gamma}{\lambda} \left(V_{\tau}^{\pi_{\tau}^*}(\rho) - \min_{n=0,1,\ldots,N-1} V_{\tau}^{\pi^N}(\rho) \right) + \frac{\alpha}{\lambda} + \frac{4}{(1-\gamma)\lambda} \sum_{n=0}^{N-1} \delta_n.$$

Hence

$$N\frac{1-\gamma}{\lambda}\Big(\min_{n=0,1,\ldots,N-1}V_{\tau}^{\pi^N}(\rho)-V_{\tau}^{\pi^*_{\tau}}(\rho)\Big)\leq \frac{\alpha}{\lambda}+y^0+\frac{4}{(1-\gamma)\lambda}\sum_{n=0}^{N-1}\delta_n.$$

and so

$$0 \leq \min_{\rho = 0, 1, \dots, N-1} V_{\tau}^{\pi^N}(\rho) - V_{\tau}^{\pi^*_{\tau}}(\rho) \leq \frac{1}{N} \frac{\alpha + \lambda y^0}{1 - \gamma} + \frac{1}{N} \frac{4}{(1 - \gamma)^2} \sum_{n=1}^{N-1} \delta_n.$$

Convergence of mirror descent with approximate advantage

Theorem 13

Given $\pi_0 \in \Pi_{\mu}$, let $(\pi_n)_{\mathbb{N}}$ be given by

$$\pi^{n+1}(\mathit{da}|s) = \operatorname*{argmin}_{m \in \mathcal{P}(A)} \int_A \Big(\hat{A}_n(s,a) + au \ln rac{\mathrm{d}\pi^n}{\mathrm{d}\mu}(a|s)\Big) (m(\mathit{da}) - \pi^n(\mathit{da}|s)) + \lambda \operatorname{\mathsf{KL}}(m|\pi^n(\cdot|s) \,.$$

where
$$\hat{A}_n(s,a) = A^{\pi_n}_{\tau}(s,a) + \mathcal{E}_n(s,a)$$
 and $\|\mathcal{E}\|_{B_b(S \times A)} = \delta_n < \infty$ for all $n \in \mathbb{N}$. Then

$$0 \leq \min_{n=0,1,\ldots,N-1} V_{\tau}^{\pi^{N}}(\rho) - V_{\tau}^{\pi^{*}}(\rho) \leq \frac{1}{N} \frac{\alpha + \lambda y^{0}}{1 - \gamma} + \frac{1}{N} \frac{4}{(1 - \gamma)^{2}} \sum_{n=0}^{N-1} \delta_{n},$$

where
$$\alpha := -\int_{\mathcal{S}} (V_{\tau}^* - V^0)(s) d_{\rho}^{\pi_{\tau}^*}(ds)$$
 and $y^0 := \int_{\mathcal{S}} \mathsf{KL}(\pi_{\tau}^* | \pi^0)(s) \ d_{\rho}^{\pi_{\tau}^*}(ds)$.

This is a small extension of results in [Kerimkulov et al., 2025a], [Lan, 2023].

Natural policy gradient is mirror descent



Natural policy gradient (NPG)

Let $\frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s) := \frac{e^{g_{\theta}(s,a)}}{Z_{\theta}(s)}$, $Z_{\theta}(s) := \int_{A} e^{g_{\theta}(s,a')}\mu(da')$ with $g_{\theta}(s,a) = (\theta,\phi(s,a))_{\mathbb{H}}$.

Fisher information matrix

$$F(heta) := \int_{\mathcal{S}} \int_{\mathcal{A}}
abla_{ heta} \ln rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu} \otimes
abla_{ heta} \ln rac{\mathrm{d}\pi_{ heta}}{\mathrm{d}\mu} (\mathsf{a}|\mathsf{s})\pi_{ heta}(\mathsf{d}\mathsf{a}|\mathsf{s}) d_{
ho}^{\pi_{ heta}}(\mathsf{d}\mathsf{a}|\mathsf{s}) \,,$$

where for $\theta, \theta' \in \mathbb{H}$ we have $(\theta \otimes \theta')_{jk} = \theta_j \theta'_k$. Let

$$\phi_{\pi_{ heta}} := \phi(s, \mathsf{a}) - \int_{\mathsf{A}} \phi(s, \mathsf{a}') \pi_{ heta}(\mathsf{d}\mathsf{a}'|s)$$
 .

Recalling (15) we have that $\nabla_{\theta} \ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s) = \nabla_{\theta} g_{\theta}(s,a) - \int_{\mathcal{A}} (\nabla_{\theta} g_{\theta})(s,a') \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a'|s)\mu(da') = \phi_{\pi_{\theta}}(s,a)$. Hence

$$F(\theta) = \int_{\mathcal{S}} \int_{\mathcal{A}} \phi_{\pi_{\theta}} \otimes \phi_{\pi_{\theta}}(s, a) \pi_{\theta}(da|s) d_{\rho}^{\pi_{\theta}}(da|s)$$
.

Natural policy gradient (NPG) updates are

$$\theta_{n+1} = \theta_n - \eta F(\theta)^{\dagger} \nabla_{\theta} V_{\tau}^{\pi_{\theta} n}(\rho), \quad n = 0, 1, \dots, \quad \theta^0 \in \mathbb{H} \quad \text{given.}$$
 (28)

Here, for $M \in \mathcal{L}(\mathbb{H}, \mathbb{H})$ we use M^{\dagger} to denote the Moore-Penrose pseudo-inverse (which coincides with M^{-1} for invertible M).

NPG in RL is due to [Kakade, 2001].

NPG is Mirror descent

Proposition 6

If given $\theta \in \mathbb{H}$ we take $\ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}(a|s) = (\theta, \phi_{\theta})_{\mathbb{H}}$ and thus obtain π_{θ_n} corresponding to θ_n then $\pi_{\theta_{n+1}}$ with θ_{n+1} given by the NPG update (28) is equal to π^{n+1} given by

$$\pi_{ heta_{n+1}}(\cdot|s) = \operatorname*{argmin}_{m \in \mathcal{P}(A)} \int_A \left(\hat{w}(heta_n) + au heta_n, \phi_{\pi_{ heta_n}}(s, a)
ight)_{\mathbb{H}} (m(da) - \pi_{ heta_n}(da|s)) + \lambda \, \mathsf{KL}(m|\pi_{ heta_n}(\cdot|s))$$

which is the mirror descent update (17) where the flat derivative is replaced by its approximation $\hat{A}_n = (\hat{w}(\theta) + \tau \theta, \phi_{\pi_0})_{\mathbb{H}}$.

Remark: Let²

$$L^{\pi_{ heta}}(w) := rac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} |A^{\pi_{ heta}}_{ au}(s, \mathsf{a}) - (w, \phi_{\pi_{ heta}}(s, \mathsf{a}))_{\mathbb{H}}|^2 \pi_{ heta}(\mathsf{d}\mathsf{a}|s) d^{\pi_{ heta}}_{
ho}(\mathsf{d}s) \,,$$

(29)

where $A_{\tau}^{\pi_{\theta}}(s,a) = Q_{\tau}^{\pi_{\theta}}(s,a) - V_{\tau}^{\pi_{\theta}}(s)$. So NPG updates are:

$$\theta_{n+1} = \theta_n, -\frac{1}{2}(\hat{w}(\theta_n) + \tau \theta_n),$$

where $\hat{w}(\theta_n)$ is the minimizer for (29).

²We are not including the $\ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu}$ term. It's just an additive term we can trivially see that $|\ln \frac{\mathrm{d}\pi_{\theta}}{\mathrm{d}\mu} - (y,\phi_{\pi_{\theta}})_{\mathbb{H}}|^2$ is minimized by $y = \theta$.

⁽²³⁾

Proof. Notice that

$$\nabla_w L^{\pi_{\theta}}(w) = \int_{\mathcal{S}} \int_{\mathcal{A}} (\mathcal{A}^{\pi_{\theta}}_{\tau}(s,a) - (w,\phi_{\pi_{\theta}}(s,a))_{\mathbb{H}}) \phi_{\pi_{\theta}}(s,a) \pi_{\theta}(da|s) d^{\pi_{\theta}}_{\rho}(ds)$$

and so the first order condition for any minimizer \hat{w} of (29) is

$$\int_{\mathcal{S}}\int_{A}(\hat{w},\phi_{\pi_{\theta}}(s,a))_{\mathbb{H}}\phi_{\pi_{\theta}}(s,a)\pi_{\theta}(\mathsf{d} \mathsf{a} | s)d_{\rho}^{\pi_{\theta}}(\mathsf{d} \mathsf{s})=\int_{\mathcal{S}}\int_{A}A_{\tau}^{\pi_{\theta}}(s,a)\phi_{\pi_{\theta}}(s,a)\pi_{\theta}(\mathsf{d} \mathsf{a} | s)d_{\rho}^{\pi_{\theta}}(\mathsf{d} \mathsf{s})\,.$$

Moreover, for any $w \in \mathbb{H}$ we have $F(\theta)w = \int_S \int_A (w, \phi_{\pi_\theta}(s, a))_{\mathbb{H}} \phi_{\pi_\theta}(s, a) \pi_\theta(da|s) d_\rho^{\pi_\theta}(ds)$. Noting also that the minimizer above depends on θ we have

$$F(\theta)\hat{w}(\theta) = \int_{S} \int_{A} A_{\tau}^{\pi_{\theta}}(s,a) \phi_{\pi_{\theta}}(s,a) \pi_{\theta}(da|s) d_{\rho}^{\pi_{\theta}}(ds).$$

Note that the Moore-Penrose pseudo-inverse provides the smallest norm solution to this i.e.

$$\hat{w}(\theta) = F(\theta)^{\dagger} \int_{\mathcal{S}} \int_{\mathcal{A}} A_{\tau}^{\pi_{\theta}}(s, \mathbf{a}) \phi_{\pi_{\theta}}(s, \mathbf{a}) \pi_{\theta}(d\mathbf{a}|s) d_{\rho}^{\pi_{\theta}}(ds).$$

This, together with (16) leads to

$$\begin{split} F(\theta)^{\dagger} \nabla_{\theta} V_{\tau}^{\pi_{\theta}}(\rho) &= \frac{1}{1 - \gamma} F(\theta)^{\dagger} \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)}^{s \sim d_{\rho}^{\pi_{\theta}}} \Big[\Big(A_{\tau}^{\pi_{\theta}}(s, a) + \tau \ln \frac{\mathrm{d} \pi_{\theta}}{\mathrm{d} \mu}(a \mid s) \Big) \phi_{\pi_{\theta}}(s, a) \Big] \\ &= \frac{1}{1 - \gamma} \big(\hat{w}(\theta) + \tau \theta \big) \,. \end{split}$$

Have

$$F(heta)^\dagger
abla_ heta V_ au^{\pi_ heta}(
ho) = rac{1}{1-\gamma} ig(\hat{w}(heta) + au hetaig)$$
 .

So the NPG stepping scheme (28) becomes

$$\theta_{n+1} = \theta_n - \frac{\eta}{1-\alpha} (\hat{w}(\theta_n) + \tau \theta_n), \quad n = 0, 1, \dots, \quad \theta_0 \in \mathbb{H}$$
 given.

Letting $\lambda = \eta (1 - \gamma)^{-1}$ we have

$$(heta_{n+1},\phi)_{\mathbb{H}}=(heta_n,\phi)_{\mathbb{H}}-rac{1}{\lambda}ig(\hat{w}(heta_n)+ au heta_n,\phi(s,a)ig)_{\mathbb{H}}\,.$$

Since $\ln \frac{\mathrm{d}\pi_{\theta_n}}{\mathrm{d}\mu}(a|s) = (\theta_n,\phi)_{\mathbb{H}} - \left(\theta_n,\int_A\phi(\cdot,a')\pi_{\theta_n}(da'|\cdot)\right)_{\mathbb{H}}$ and collecting all the terms constant in a in some b=b(s)

we then have

$$\ln \frac{\mathrm{d}\pi_{\theta_{n+1}}}{\mathrm{d}\mu}(a|s) = \ln \frac{\mathrm{d}\pi_{\theta_{n}}}{\mathrm{d}\mu}(a|s) - \frac{1}{\lambda} \left(\hat{w}(\theta_{n}) + \tau\theta_{n}, \phi_{\pi_{\theta_{n}}}(s, a)\right)_{\mathbb{H}} + b(s),$$

with b chosen such that $\pi_{\theta_{n+1}} \in \mathcal{P}(A|S)$. Hence

$$\ln rac{\mathrm{d}\pi_{ heta_n+1}}{\mathrm{d}\pi_{\pi_{ heta_n}}}(a|s) = -rac{1}{\lambda}ig(\hat{w}(heta_n) + au heta_n, \phi_{\pi_{ heta_n}}(s,a)ig)_{\mathbb{H}} + b(s)\,.$$

And so

$$rac{\mathrm{d}\pi_{ heta_{n+1}}}{\mathrm{d}\pi_{\pi_{ heta}}}(a|s) = \exp\left(-rac{1}{\lambda}\left(\hat{w}(heta_n) + au heta_n, \phi_{\pi_{ heta_n}}(s,a)
ight)_{\mathbb{H}} + b(s)
ight).$$

Then

$$\pi_{\theta_{n+1}}(\cdot|s) = \operatorname*{argmin}_{m \in \mathcal{D}(A)} \int_{A} \left(\hat{w}(\theta_n) + \tau \theta_n, \phi_{\pi_{\theta_n}}(s, a) \right)_{\mathbb{H}} (m(da) - \pi_{\theta_n}(da|s)) + \lambda \operatorname{\mathsf{KL}}(m|\pi_{\theta_n}(\cdot|s)),$$

due to [Dupuis and Ellis, 1997], Lemma 1.4.3. □

PPO and FR-PPO



Connection to PPO

Proximal policy optimization (PPO) [Schulman et al., 2017] optimizes:

$$J_{ ext{PPO}}(heta) = \mathbb{E}_{a \sim \pi_{ heta_{ ext{old}}}(\cdot | s)}^{s \sim \pi_{ heta_{ ext{old}}}(\cdot | s)} igg(\min \left[rac{\delta V_{ au}^{\pi_{ heta_{ ext{old}}}}(s, a) rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \pi_{ heta_{ ext{old}}}}(a | s),
ight. \ \left. \left. \operatorname{clip} \left(1 - arepsilon, 1 + arepsilon, rac{\mathrm{d} \pi_{ heta}}{\mathrm{d} \pi_{ heta_{ ext{old}}}}(a | s)
ight) rac{\delta V_{ au}^{\pi_{ heta_{ ext{old}}}}(s, a)}{\delta \pi} igg)
ight]
ight).$$

Proximal policy optimization algorithms
J.Schulman. F. Wolski. P. Dhariwal. A. Radford... - arXiv r
... call proximal policy optimization (PPO), have some
Our experiments test PPO on a collection of benchmark
\$\tilde{\Sigma}\$ Save \$90 Cite Cited by \$\frac{25072}{25072}\$ Related articles

DeepSeek-AI

research@deepseek.com

DeepSeek-V3 Technical Report

$$\frac{1}{G} \sum_{i=1}^{G} \left(\min \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)} A_i, \operatorname{clip} \left(\frac{\pi_{\theta}(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) \right)$$

TRPO

- Switch to maximizing in line with [Lascu et al., 2025].
- $A_{\pi}=Q_{\pi}-V_{\pi}=rac{\delta V^{\pi}}{\delta \pi}$, (objective has au=0 i.e. no KL term).

Theorem 14 ([Achiam et al., 2017])

For all $\rho \in \mathcal{P}(S)$ and any policies $\pi', \pi \in \mathcal{P}(A)$,

$$V^{\pi'}(
ho) - V^{\pi}(
ho) \geq rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{A} rac{\mathrm{d}\pi'}{\mathrm{d}\pi} A_{\pi}(s,a) \pi(\mathrm{d}a|s) d^{\pi}_{
ho}(\mathrm{d}s) \ - rac{4\gamma \|r\|_{\mathcal{B}_{b}(S imes A)}}{(1-\gamma)^{3}} \int_{\mathcal{S}} \mathsf{TV}(\pi'(\cdot|s), \pi(\cdot|s)) d^{\pi}_{
ho}(\mathrm{d}s) =: J(\pi').$$

N.B.
$$J(\pi) = 0$$
 so $\max_{\pi'} J(\pi') > 0$.

So

$$\pi' \in \operatorname*{argmax} J(\pi')$$

guarantees improvement: $V^{\pi'}-V^{\pi}=\max_{\pi'}J(\pi')\geq 0.$

N.B.

$$\int_{\mathcal{S}} \mathsf{TV}(\pi'(\cdot|s), \pi(\cdot|s)) d_{\rho}^{\pi}(\mathrm{d}s) = \tfrac{1}{2} \int_{\mathcal{S}} \int_{\mathcal{A}} \left| \tfrac{\mathrm{d}\pi'}{\mathrm{d}\pi}(\mathsf{a}|s) - 1 \right| \pi(\mathrm{d}\mathsf{a}|s) d_{\rho}^{\pi}(\mathrm{d}s) \,,$$

we can write surrogate loss as

$$V^{\pi'}(
ho) - V^{\pi}(
ho) \geq rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{A} rac{\mathrm{d}\pi'}{\mathrm{d}\pi} A_{\pi}(s,a) \pi(\mathrm{d}a|s) d^{\pi}_{
ho}(\mathrm{d}s) \ - rac{2\gamma \|r\|_{B_{b}(\mathcal{S} imes A)}}{(1-\gamma)^{3}} \int_{\mathcal{S}} \int_{A} \left| rac{\mathrm{d}\pi'}{\mathrm{d}\pi} (a|s) - 1 \right| \pi(\mathrm{d}a|s) = J(\pi').$$

This is a motivation for PPO's constraining the probability ratio $\left|\frac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s)-1\right|\leq \varepsilon$, where $\varepsilon>0$ is a hyperparameter.

PPO

$$J_{\mathsf{PPO}}(\pi') \coloneqq \int_{S} \int_{A} \min\left\{ rac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s) A_{\pi}(s,a), \mathsf{clip}_{1-arepsilon}^{1+arepsilon}\left(rac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s)
ight) A_{\pi}(s,a)
ight\} \pi(\mathrm{d}a|s) d_{
ho}^{\pi}(\mathrm{d}s) \,.$$

The challenge of TV geometry

$$J(\pi') = \frac{1}{1-\gamma} \int_{\mathcal{S}} \int_{\mathcal{A}} \frac{\mathrm{d}\pi'}{\mathrm{d}\pi} A_{\pi}(s,a) \pi(\mathrm{d}a|s) d_{\rho}^{\pi}(\mathrm{d}s) - \frac{2\gamma \|r\|_{\mathcal{B}_{b}(S\times\mathcal{A})}}{(1-\gamma)^{3}} \int_{\mathcal{S}} \int_{\mathcal{A}} \left| \frac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s) - 1 \right| \pi(\mathrm{d}a|s).$$

Seen how to get convergence if penalty satisfies 3-point property (Bregman proximal inequality).

Another approach is to go via first order optimality condition. But ...

$$2\frac{\delta \operatorname{TV}(\pi',\pi)(s)}{\delta \pi'}(a) = \frac{\delta}{\delta \pi'} \left[\int_{S} \int_{A} \left| \frac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s) - 1 \right| \pi(\mathrm{d}a|s) \right] = \operatorname{sign}\left(\frac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s) - 1 \right) \frac{\mathrm{d}\pi}{\mathrm{d}\lambda}(a|s).$$

Therefore, the first-order condition for the update rule gives the scheme

$$A_\pi(s,a) - rac{1}{ au}\operatorname{sign}\left(rac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s) - 1
ight)rac{\mathrm{d}\pi}{\mathrm{d}\lambda}(a|s) = ext{const in a}.$$

ullet Two solutions for $rac{\mathrm{d}\pi'}{\mathrm{d}\pi}$ depending on whether $rac{\mathrm{d}\pi'}{\mathrm{d}\pi}(a|s)-1\geq 0$ or not.

Improved lower bound

Theorem 15 (Lower bound on performance difference [Lascu et al., 2025])

For any π' , $\pi \in \mathcal{P}_{\lambda}(A|S)$, it holds that

$$V^{\pi'}(
ho) - V^{\pi}(
ho) \geq rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{\mathcal{A}} rac{\mathrm{d}\pi'}{\mathrm{d}\pi} (\mathsf{a}|\mathsf{s}) \mathcal{A}_{\pi}(\mathsf{s},\mathsf{a}) \pi(\mathrm{d}\mathsf{a}|\mathsf{s}) \mathcal{d}^{\pi}_{
ho}(\mathrm{d}\mathsf{s}) \ - rac{8\gamma \|r\|_{\mathcal{B}_{b}(S imes A)}}{(1-\gamma)^{3}} \int_{\mathcal{S}} \mathsf{TV}^{2}(\pi'(\cdot|\mathsf{s}),\pi(\cdot|\mathsf{s})) \mathcal{d}^{\pi}_{
ho}(\mathrm{d}\mathsf{s}).$$

Towards Fisher-Rao PPO

Fix $\lambda \in \mathcal{M}_+(A)$. Then the squared Fisher-Rao (FR) distance between measures μ and μ' is defined by

$$\mathsf{FR}^2(\mu,\mu') = 4 \int_{\mathcal{A}} \left| \sqrt{rac{\mathrm{d}\mu}{\mathrm{d}\lambda}} - \sqrt{rac{\mathrm{d}\mu'}{\mathrm{d}\lambda}} \right|^2 \mathrm{d}\lambda\,,$$

if $\mu, \mu' \ll \lambda$, and ∞ otherwise.

Cauchy-Schwarz inequality gives

$$\int_S \mathsf{TV}^2(\pi'(\cdot|s),\pi(\cdot|s)) d_\rho^\pi(\mathrm{d} s) \leq \frac{1}{16} \int_S \mathsf{FR}^2(\pi'(\cdot|s)^2,\pi(\cdot|s)^2) d_\rho^\pi(\mathrm{d} s),$$

Corollary 16

For any π' , $\pi \in \mathcal{P}_{\lambda}(A|S)$, it holds that

$$egin{aligned} V^{\pi'}(
ho) - V^{\pi}(
ho) &\geq rac{1}{1-\gamma} \int_{\mathcal{S}} \int_{\mathcal{A}} rac{\mathrm{d}\pi'}{\mathrm{d}\pi} (\mathsf{a}|\mathsf{s}) \mathcal{A}_{\pi}(\mathsf{s},\mathsf{a}) \pi(\mathrm{d}\mathsf{a}|\mathsf{s}) d^{\pi}_{
ho}(\mathrm{d}\mathsf{s}) \ &- rac{\gamma \|r\|_{\mathcal{B}_{b}(\mathcal{S} imes\mathcal{A})}}{2(1-\gamma)^{3}} \int_{\mathcal{S}} \mathsf{FR}^{2} (\pi'(\cdot|\mathsf{s})^{2},\pi(\cdot|\mathsf{s})^{2}) d^{\pi}_{
ho}(\mathrm{d}\mathsf{s}), \end{aligned}$$

where $\pi(\cdot|s)^2 := \left(\frac{\mathrm{d}\pi}{\mathrm{d}\lambda}\right)^2(\cdot|s)$, for all $s \in S$.

Fisher-Rao PPO

Consider

$$\pi^{n+1}(\cdot|s) = \operatorname*{argmax}_{m \in \mathcal{P}} \left[\int_{A} A_{\pi^n}(s,a) \frac{\mathrm{d} m}{\mathrm{d} \pi^n}(a|s) \pi^n(da|s) - \frac{1}{2\tau} \operatorname{FR}^2(m^2,\pi^n(\cdot|s)^2) \right].$$

Theorem 17 (Policy improvement [Lascu et al., 2025])

Let $V^n:=V^{\pi^n}$. If $\frac{1}{\tau}\geq \frac{\|r\|_{B_b(S imes A)}}{(1-\gamma)^2}$, then for any $\rho\in\mathcal{P}(S)$ we have

$$V^{n+1}(\rho) \geq V^n(\rho)$$
 for all $n > 0$.

Theorem 18 (Sublinear convergence [Lascu et al., 2025])

Let $V^n:=V^{\pi^n}$. For all $\rho\in\mathcal{P}(S)$, $\frac{1}{\tau}\geq \frac{\gamma\|r\|_{B_b(S imes A)}}{(1-\gamma)^2}$, for any $\pi\in\mathcal{P}(A|S)$ and $n\in\mathbb{N}_0$,

$$(V^\pi - V^n)(
ho) \leq rac{1}{n(1-\gamma)} igg(rac{1}{ au} \int_S \mathsf{FR}^2((\pi(\cdot|s)^2, \pi_0(\cdot|s)^2) d^\pi_
ho(ds) + \int_S (V^0 - V^\pi)(s) d^\pi_
ho(ds)igg)\,.$$

Summary

General

- MDPs can be formulated in very general (Polish) state and action spaces.
- Relaxed MDP with entropy regularization has good mathematical properties (existence and uniqueness of optimal policy) under minimal assumptions (boundedness and measurability).

Policy gradient

- We have PG theorem but convergence analysis is hard due to non-convexity (best known results need non-local Lojasiewicz are in finite state & action spaces or structural assumptions - LQR).
- Mirror descent (NPG) is much more ameaneable to convergence analysis (with direct parametrization of linear basis functions for log gensities), key ingredients
 - Three point lemma
 - L-smoothness (from perf. diff. and Lan's trick)
 - The "replacement-for-convexity-property" (that perf. diff provides)
- Convergence of popular PG schemes e.g. PPO still not understood.
- Replacing linear parametrizations (basis functions, NTK, mean-field) with deep networks is ... not understood even for supervised learning.

Thank you!



References I

- [Achiam et al., 2017] Achiam, J., Held, D., Tamar, A., and Abbeel, P. (2017). Constrained policy optimization. In *International conference on machine learning*, pages 22–31. PMLR.
- [Agarwal et al., 2019] Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2019). Optimality and approximation with policy gradient methods in Markov decision processes. arXiv preprint arXiv:1908.00261.
- [Aubin-Frankowski et al., 2022] Aubin-Frankowski, P.-C., Korba, A., and Léger, F. (2022). Mirror descent with relative smoothness in measure spaces, with application to Sinkhorn and EM. *Advances in Neural Information Processing Systems*, 35:17263–17275.
- [Beck and Teboulle, 2003] Beck, A. and Teboulle, M. (2003). Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175.
- [Bu et al., 2019] Bu, J., Mesbahi, A., Fazel, M., and Mesbahi, M. (2019). LQR through the lens of first order methods: Discrete-time case. arXiv preprint arXiv:1907.08921.
- [Cayci et al., 2021] Cayci, S., He, N., and Srikant, R. (2021). Linear convergence of entropy-regularized natural policy gradient with linear function approximation. arXiv preprint arXiv:2106.04096.
- [Cen et al., 2022] Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2022). Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578.
- [Doya, 2000] Doya, K. (2000). Reinforcement learning in continuous time and space. Neural computation, 12(1):219–245.

References II

- [Dupuis and Ellis, 1997] Dupuis, P. and Ellis, R. S. (1997). A weak convergence approach to the theory of large deviations. John Wiley & Sons, Inc., New York.
- [Fazel et al., 2018] Fazel, M., Ge, R., Kakade, S., and Mesbahi, M. (2018). Global convergence of policy gradient methods for the linear quadratic regulator. In *International Conference on Machine Learning*, pages 1467–1476. PMLR.
- [Geist et al., 2019] Geist, M., Scherrer, B., and Pietquin, O. (2019). A theory of regularized Markov decision processes. In International Conference on Machine Learning, pages 2160–2169. PMLR.
- [Giegrich et al., 2024] Giegrich, M., Reisinger, C., and Zhang, Y. (2024). Convergence of policy gradient methods for finite-horizon exploratory linear-quadratic control problems. SIAM Journal on Control and Optimization, 62(2):1060–1092.
- [Haarnoja et al., 2017] Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. (2017). Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pages 1352–1361. PMLR.
- [Haarnoja et al., 2018] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR.
- [Hernández-Lerma and Lasserre, 2012] Hernández-Lerma, O. and Lasserre, J. B. (2012). Discrete-time Markov control processes: basic optimality criteria, volume 30. Springer Science & Business Media.
- [Howard, 1960] Howard, R. A. (1960). Dynamic programming and markov processes. John Wiley.

References III

- [Hu et al., 2023] Hu, B., Zhang, K., Li, N., Mesbahi, M., Fazel, M., and Başar, T. (2023). Toward a theoretical foundation of policy optimization for learning control policies. *Annual Review of Control, Robotics, and Autonomous Systems*, 6:123–158.
- [Kakade and Langford, 2002] Kakade, S. and Langford, J. (2002). Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 267–274.
- [Kakade, 2001] Kakade, S. M. (2001). A natural policy gradient. Advances in neural information processing systems, 14.
- [Kerimkulov et al., 2025a] Kerimkulov, B., Leahy, J.-M., Šiška, D., Szpruch, L., and Zhang, Y. (2025a). A Fisher–Rao gradient flow for entropy-regularised Markov decision processes in Polish spaces. Foundations of Computational Mathematics.
- [Kerimkulov et al., 2025b] Kerimkulov, B., Šiška, D., Szpruch, Ł., and Zhang, Y. (2025b). Mirror descent for stochastic control problems with measure-valued controls. Stochastic Processes and their Applications, page 104765.
- [Khodadadian et al., 2022] Khodadadian, S., Jhunjhunwala, P. R., Varma, S. M., and Maguluri, S. T. (2022). On linear and super-linear convergence of natural policy gradient algorithm. Systems & Control Letters, 164:105214.
- [Lan, 2023] Lan, G. (2023). Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, 198(1):1059–1106.
- [Lascu et al., 2025] Lascu, R.-A., Šiška, D., and Szpruch, Ł. (2025). PPO in the Fisher–Rao geometry. arXiv preprint arXiv:2506.03757.

References IV

- [Lemaréchal, 2012] Lemaréchal, C. (2012). Cauchy and the gradient method. Doc Math Extra, 251(254):10.
- [Manna et al., 2022] Manna, S., Loeffler, T. D., Batra, R., Banik, S., Chan, H., Varughese, B., Sasikumar, K., Sternberg, M., Peterka, T., Cherukara, M. J., et al. (2022). Learning in continuous action space for developing high dimensional potential energy models. *Nature communications*, 13(1):368.
- [Mei et al., 2021] Mei, J., Gao, Y., Dai, B., Szepesvari, C., and Schuurmans, D. (2021). Leveraging non-uniformity in first-order non-convex optimization. In International Conference on Machine Learning, pages 7555–7564. PMLR.
- [Mei et al., 2020] Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. (2020). On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829. PMLR.
- [Nemirovski, 1979] Nemirovski, A. (1979). Efficient methods. Ekonomika i Mat. Metody, 15.
- [Schulman et al., 2015] Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. (2015). High-dimensional continuous control using generalized advantage estimation. arXiv preprint arXiv:1506.02438.
- [Schulman et al., 2017] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.
- [Sutton and Barto, 2018] Sutton, R. S. and Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.
- [Sutton et al., 1999] Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12.

References V

- [Tomar et al., 2020] Tomar, M., Shani, L., Efroni, Y., and Ghavamzadeh, M. (2020). Mirror descent policy optimization. arXiv preprint arXiv:2005.09814.
- [Van Hasselt, 2012] Van Hasselt, H. (2012). Reinforcement learning in continuous state and action spaces. In Reinforcement Learning: State-of-the-Art, pages 207–251. Springer.
- [Xiao, 2022] Xiao, L. (2022). On the convergence rates of policy gradient methods. arXiv preprint arXiv:2201.07443.
- [Zhan et al., 2023] Zhan, W., Cen, S., Huang, B., Chen, Y., Lee, J. D., and Chi, Y. (2023). Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. SIAM Journal on Optimization, 33(2):1061–1091.

Q-learning

DPP equation for Q_{τ}^* :

$$Q_{\tau}^*(s,a) = c(s,a) + \gamma \int_{S} \inf_{m \in \mathcal{P}(A)} \left(Q_{\tau}^*(s',a') + \tau \frac{\mathrm{d}m}{\mathrm{d}\mu}(a') \right) m(da') \right) P(ds'|s,a).$$

Re-write:

$$0 = \mathbb{E}_{a' \sim \pi_{\tau}^*(\cdot|s')}^{s' \sim P(\cdot|s,a)} \left[c(s,a) + \gamma \left(Q_{\tau}^*(s',a') + \tau \frac{\mathrm{d}\pi_{\tau}^*}{\mathrm{d}\mu} (a'|s') \right) \right) - Q_{\tau}^*(s,a) \right].$$

Q-learning:

$$Q_{k+1}(s_k,a_k) = Q_k(s_k,a_k) + \delta_k \left[c(s_k,a_k) + \gamma \left(Q_k(s_{k+1},a_{k+1}) + \tau \frac{\mathrm{d}\pi_k}{\mathrm{d}\mu} (a_{k+1}|s_{k+1}) \right) - Q_k(s_k,a_k) \right],$$

where $s_{k+1} \sim P(\cdot|s_k, a_k)$, $\pi_k(da'|s_k) \propto \exp(-\tau^{-1}Q_k(s_k, a'))\mu(da')$, $a_{k+1} \sim \pi_k(\cdot|s_k)$.

 $\textbf{Convergence:} \ \ \text{like value iteration} \ + \ \text{stochastic approximation (Robbins-Monro)}.$