Statistical Foundations of Reinforcement Learning: I

COLT 2021

Akshay Krishnamurthy (M Wen Sun (Cornell

Akshay Krishnamurthy (MSR, akshaykr@microsoft.com)

Wen Sun (Cornell, ws455@cornell.edu)

Reinforcement Learning: Motivation and empirical progress



TD Gammon [Tesauro]



Stratospheric balloons [Bellemare et.al]



DeepMind Starcraft [Vinyals et.al]



OpenAl Dexterous manipulation [Akkaya et.al]

Learning Agent





Environment



Determine **action** based on **state**

Learning Agent





Send **reward** and **next state**



Environment



Learning Agent









Environment



Learning Agent









Environment

Learning Agent









Environment

Learning Agent









Environment

Learning Agent













Policy search methods;





Plan for the tutorial

Part 1: Tabular setting

- Basics and key concepts 1.
- Policy optimization and Natural Policy Gradient 2.
- UCB-Value Iteration 3.

Part 2: Problem set

Part 3: Function approximation + Exploration

- Linear methods and complexity 1.
- 2.

Nonlinear methods, bellman rank, bilinear classes, representation learning

Part 1A: MDP Basics

Markov Decision Processes (Discounted version)

Learning Agent



policy $\pi(a \mid s)$

Determine action based on state



Infinitely many steps



Send reward and next state

 $r(s,a), s' \sim P(\cdot \mid s,a)$

Environment



 $\mathcal{M} = \{S, A, P, r, \gamma, \mu\}$ $\mu \in \Delta(S)$ $P : S \times A \mapsto \Delta(S)$ $r : S \times A \rightarrow [0,1]$ $\gamma \in [0,1)$

Markov Decision Processes (Discounted version)

Learning Agent



policy $\pi(a \mid s)$

Determine action based on state



Infinitely many steps



Send reward and next state

 $r(s, a), s' \sim P(\cdot \mid s, a)$



Environment



 $\mathscr{M} = \{S, A, P, r, \gamma, \mu\}$ $\mu \in \Delta(S)$ $P: S \times A \mapsto \Delta(S)$ $r: S \times A \rightarrow [0,1]$ $\gamma \in [0,1)$

Objective: $\max_{\pi} \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h}) \, | \, s_{0} \sim \mu, a_{h} \sim \pi(\, . \, | \, s_{h}), s_{h+1} \sim P(\, . \, | \, s_{h}, a_{h}) \right]$



Average State-action Distributions

Given a policy $\pi: S \mapsto \Delta(A)$

Denote $d^{\pi}_{\mu,h}(s,a) := P^{\pi}((s_h,a_h) = (s,a))$, i.e., probability of π hitting (s,a) at time step h

Average State-action Distributions

Denote
$$d^{\pi}_{\mu,h}(s,a) := P^{\pi}((s_h,a_h) = (s,a)$$

Denote
$$d^{\pi}_{\mu}(s, a) := (1 - \gamma) \sum_{h=0}^{\infty} \gamma^{h} d^{\pi}_{h}$$

- Given a policy $\pi : S \mapsto \Delta(A)$
 -), i.e., probability of π hitting (s, a) at time step h
 - ${}_{h}^{\pi}(s, a)$ as the average state-action distribution

Average State-action Distributions

Denote
$$d_{\mu,h}^{\pi}(s,a) := P^{\pi}\left((s_h,a_h) = (s,a)\right)$$

Denote $d_{\mu}^{\pi}(s,a) := (1-\gamma) \sum_{h=0}^{\infty} \gamma^h d_h^{\pi}$

- Given a policy $\pi : S \mapsto \Delta(A)$
 -), i.e., probability of π hitting (s, a) at time step h
 - $\frac{\pi}{s,a}$ as the average state-action distribution

We will abuse notation a bit and denote $d^{\pi}_{\mu}(s) := \sum d^{\pi}_{\mu}(s, a)$ as the average state-distribution \mathcal{A}

Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

$$V^{\pi}(s) = \mathbb{E}\left[\left|\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h})\right| s_{0} = s, a_{h} \sim \pi(s_{h}), s_{h+1} \sim P(\cdot \mid s_{h}, a_{h})\right]$$

Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

$$V^{\pi}(s) = \mathbb{E}\left[\left|\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h})\right| s_{0} = s, a_{h} \sim \pi(s_{h}), s_{h+1} \sim P(\cdot \mid s_{h}, a_{h})\right]$$
$$= \mathbb{E}_{a \sim \pi(\cdot \mid s)}\left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)}V^{\pi}(s')\right]$$
(Belly)

man equation)



Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

$$V^{\pi}(s) = \mathbb{E}\left[\left|\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h})\right| s_{0} = s, a_{h} \sim \pi(s_{h}), s_{h+1} \sim P(\cdot \mid s_{h}, a_{h})\right]$$
$$= \mathbb{E}_{a \sim \pi(\cdot \mid s)}\left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^{\pi}(s')\right]$$
(Belly)

Q function $Q^{\pi}(s, a)$: total reward when starting in state s and action a and following π afterwards

man equation)



Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

$$V^{\pi}(s) = \mathbb{E}\left[\left|\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h})\right| s_{0} = s, a_{h} \sim \pi(s_{h}), s_{h+1} \sim P(\cdot \mid s_{h}, a_{h})\right]$$
$$= \mathbb{E}_{a \sim \pi(\cdot \mid s)}\left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^{\pi}(s')\right]$$
(Bella

 $Q^{\pi}(s,a) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^{h} r(s_{h},a_{h})\right](s_{0},a_{h})$

man equation)

Q function $Q^{\pi}(s, a)$: total reward when starting in state s and action a and following π afterwards

$$a_0 = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot | s_h, a_h)$$



Value function $V^{\pi}(s)$: total reward when starting in state s and following π afterwards

$$V^{\pi}(s) = \mathbb{E}\left[\left|\sum_{h=0}^{\infty} \gamma^{h} r(s_{h}, a_{h})\right| s_{0} = s, a_{h} \sim \pi(s_{h}), s_{h+1} \sim P(\cdot \mid s_{h}, a_{h})\right]$$
$$= \mathbb{E}_{a \sim \pi(\cdot \mid s)}\left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^{\pi}(s')\right]$$
(Bella

 $Q^{\pi}(s,a) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^{h} r(s_{h},a_{h})\right](s_{0},a_{h})$

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

man equation)

Q function $Q^{\pi}(s, a)$: total reward when starting in state s and action a and following π afterwards

$$a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot | s_h, a_h)$$

(Bellman equation)



There exists a deterministic stationary policy $\pi^* : S \mapsto A$, s.t., $V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$

Optimality

We denote $V^{\star} := V^{\pi^{\star}}$

Optimality

There exists a deterministic stationary policy $\pi^* : S \mapsto A$, s.t., $V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$

$$:= V^{\pi^*}, Q^* := Q^{\pi^*}$$

We denote V^{\star}

Theorem 1: Bellman Optimality

 $\forall s, a : Q^{\star}(s, a) = r(s, a)$

Optimality

There exists a deterministic stationary policy $\pi^* : S \mapsto A$, s.t., $V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$

$$:= V^{\pi^*}, Q^* := Q^{\pi^*}$$

$$(a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q^{\star}(s',a')$$

We denote V^{\star}

Theorem 1: Bellman Optimality

 $\forall s, a : Q^{\star}(s, a) = r(s, a)$

Theorem 2: Bellman Optimality

For any $Q: S \times A \rightarrow \mathbb{R}$, if Q(s, a) =

for all s, a, then Q(s,

Optimality

There exists a deterministic stationary policy $\pi^* : S \mapsto A$, s.t., $V^{\pi^{\star}}(s) \geq V^{\pi}(s), \forall s, \pi$

$$:= V^{\pi^*}, Q^* := Q^{\pi^*}$$

$$(a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q^{\star}(s',a')$$

$$= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \max_{a'} Q(s', a')$$
$$a) = Q^{\star}(s, a), \forall s, a$$



Planning in MDP with known transition P and reward r

i.e., how to compute π^{\star} (and V^{\star} / Q^{\star}) given the MDP (P, r)

Idea: fixed point iteration

Define: Bellman operator $\mathcal{T}: (S \times A \to \mathbb{R}) \to (S \times A)$

$(\mathcal{T}f)_{s,a} := r(s,a) + \gamma \mathbb{E}_{s'}$

$$\stackrel{\rightarrow}{} (S \times A \rightarrow \mathbb{R})$$

$$\stackrel{\prime}{}_{a'} \stackrel{}{} [\max_{a'} f(s', a')]$$

Idea: fixed point iteration

Define: Bellman operator $\mathcal{T}: (S \times A \to \mathbb{R}) (\mathcal{T}f)_{s,a} := r(s,a) + \gamma \mathbb{E}_{s'}$

VI Algorithm: Initialize $Q^{(0)}$ s.t., $Q^{(0)}(s, a) \in [0, 1/(1 - \gamma))$ Iterate $Q^{(t+1)} \leftarrow \mathcal{T}Q^{(t)}$

$$(S \times A \to \mathbb{R})$$

$$(\sim_{P(\cdot|s,a)} [\max_{a'} f(s',a')]$$

Idea: fixed point iteration

Define: Bellman operator $\mathscr{T} : (S \times A \to \mathbb{R}) - (\mathscr{T}f)_{s,a} := r(s,a) + \gamma \mathbb{E}_{s'}$

VI Algorithm: Initialize $Q^{(0)}$ s.t., $Q^{(0)}(s, a) \in [0, 1/(1 - \gamma))$ Iterate $Q^{(t+1)} \leftarrow \mathcal{T}Q^{(t)}$

> Theorem: Induced policy $\pi^{(t)} : s \mapsto \arg \max_{a} Q^{(t)}(s, a)$ satisfies $V^{\pi^{(t)}}(s) \ge V^{\star}(s) - \frac{2\gamma^{t}}{1-\gamma} \|Q^{(0)} - Q^{\star}\|_{\infty} \quad \forall s \in S$

$$(S \times A \to \mathbb{R})$$

$$\sim_{P(\cdot|s,a)} [\max_{a'} f(s',a')]$$

Idea: fixed point iteration

Define: Bellman operator $\mathscr{T} : (S \times A \to \mathbb{R}) - (\mathscr{T}f)_{s,a} := r(s,a) + \gamma \mathbb{E}_{s'}$

VI Algorithm: Initialize $Q^{(0)}$ s.t., $Q^{(0)}(s, a) \in [0]$ Iterate $Q^{(t+1)} \leftarrow \mathcal{T}Q^{(t)}$

> Theorem: Induced policy $\pi^{(t)} : s \mapsto \arg n$ $V^{\pi^{(t)}}(s) \ge V^{\star}(s) - \frac{2\gamma^t}{1-\gamma} \|Q^{(0)} - Q^{\star}\|_{\infty}$



MDP Planning: Policy iteration

Idea: Alternate between policy evaluation and policy improvement Initialize $\pi^{(0)}: S \to A$

Repeat:

- Compute $Q^{\pi^{(t)}}$ (evaluation)
- Update $\pi^{(t+1)}$: $\pi^{(t+1)}(s) = \arg \max Q^{\pi^{(t)}}(s, a)$ (improvement)
MDP Planning: Policy iteration

Idea: Alternate between policy evaluation and policy improvement Initialize $\pi^{(0)}: S \to A$

Repeat:

- Compute $Q^{\pi^{(t)}}$ (evaluation)
- Update $\pi^{(t+1)}$: $\pi^{(t+1)}(s) = \arg \max Q^{\pi^{(t)}}(s, a)$ (improvement)



MDP Planning: Policy iteration

Idea: Alternate between policy evaluation and policy improvement Initialize $\pi^{(0)}: S \to A$

Repeat:

- Compute $Q^{\pi^{(t)}}$ (evaluation)
- Update $\pi^{(t+1)}$: $\pi^{(t+1)}(s) = \arg \max Q$ \boldsymbol{a}

Theorem: Geometric convergence:

$$\|V^{\pi^{(t+1)}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^{(t)}}\|_{\infty}$$



$$P^{\pi^{(t)}}(s,a)$$
 (improvement)

$$-V^{\star}\|_{\infty}$$

Finite Horizon MDPs

- $\mathscr{M} = \{S, A, P, r, \mu, H\}$

$P: S \times A \mapsto \Delta(S), \quad r: S \times A \to [0,1], \quad H \in \mathbb{N}^+, \quad \mu \in \Delta(S)$

time-dependent policies: $\pi^* := \{\pi_0^*, \dots, \pi_{H-1}^*\}$ time-dependent V/Q functions: $\{V_h^{\star}\}_{h=0}^{H-1}, \{Q_h^{\star}\}_{h=0}^{H-1}$

Finite Horizon MDPs

- $\mathscr{M} = \{S, A, P, r, \mu, H\}$

Episode:

- $s_0 \sim \mu$ For h = 0, ..., H - 1:
- Take action a_h
- Collect reward $r(s_h, a_h)$
- Transition $s_{h+1} \sim P(\cdot \mid s_h, a_h)$

$P: S \times A \mapsto \Delta(S), \quad r: S \times A \to [0,1], \quad H \in \mathbb{N}^+, \quad \mu \in \Delta(S)$

time-dependent policies: $\pi^* := \{\pi_0^*, \dots, \pi_{H-1}^*\}$ time-dependent V/Q functions: $\{V_h^{\star}\}_{h=0}^{H-1}, \{Q_h^{\star}\}_{h=0}^{H-1}$

Finite Horizon MDPs

- $\mathscr{M} = \{S, A, P, r, \mu, H\}$
- $P: S \times A \mapsto \Delta(S), \quad r: S \times A \to [0,1], \quad H \in \mathbb{N}^+, \quad \mu \in \Delta(S)$

Episode:

- $s_0 \sim \mu$ O For h = 0, ..., H - 1:
- Take action a_h
- Collect reward $r(s_h, a_h)$
- Transition $s_{h+1} \sim P(\cdot \mid s_h, a_h)$

bjective function:
$$V(\pi) = \mathbb{E} \left[\sum_{h=0}^{H-1} r(s_h, a_h) \right]$$

time-dependent policies: $\pi^* := \{\pi_0^*, \dots, \pi_{H-1}^*\}$ time-dependent V/Q functions: $\{V_h^{\star}\}_{h=0}^{H-1}, \{Q_h^{\star}\}_{h=0}^{H-1}$

Summary so far:

- MDP definitions (discounted infinite horizon & finite horizon);
- State-action distributions, value and Q functions, and two planning algorithms

Part 1B: Policy Gradient & Natural Policy Gradient

Policy Optimization Motivation: Practical





[OpenAl Five, 18]

[OpenAl, 19]

Policy Optimization Motivation: Simple

- $\pi_{\theta}(a \mid s) := \pi(a \mid s; \theta)$
 - $\theta_{t+1} = \theta_t$

$$\begin{array}{l} \mathcal{P}) \quad V^{\pi_{\theta}} = \mathbb{E}_{\pi_{\theta}} \left[\sum_{h=0}^{\infty} \gamma^{h} r_{h} \right] \\ + \eta \nabla_{\theta} V^{\pi_{\theta}} |_{\theta = \theta_{t}} \end{array}$$

Policy Optimization Motivation: Simple

- $\pi_{\theta}(a \mid s) := \pi(a \mid s; \theta)$
 - $\theta_{t+1} = \theta_t$

We can have a closed-form expression for PG:

- Define advantage function $A^{\pi_{\theta}}(s)$

$$\nabla_{\theta} V^{\pi_{\theta}} = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d_{\mu}^{\pi_{\theta}}} \left[\nabla_{\theta} \ln \pi_{\theta}(a \mid s) A^{\pi_{\theta}}(s, a) \right]$$

$$\begin{array}{l} \mathcal{P}) \quad V^{\pi_{\theta}} = \mathbb{E}_{\pi_{\theta}} \left[\sum_{h=0}^{\infty} \gamma^{h} r_{h} \right] \\ + \eta \nabla_{\theta} V^{\pi_{\theta}} |_{\theta = \theta_{t}} \end{array}$$

Policy Gradient Theorem [Sutton, McAllester, Singh, Mansour]:

$$(s, a) := Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$$
, we have:

Policy Optimization Motivation: Simple

- $\pi_{\theta}(a \mid s) := \pi(a \mid s; \theta)$
 - $\theta_{t+1} = \theta_t$

We can have a closed-form expression for PG:



Define advantage function $A^{\pi_{\theta}}(s)$

$$\nabla_{\theta} V^{\pi_{\theta}} = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d_{\mu}^{\pi_{\theta}}} \left[\nabla_{\theta} \ln \pi_{\theta}(a \mid s) A^{\pi_{\theta}}(s, a) \right]$$

$$\begin{array}{l} \mathcal{P}) \quad V^{\pi_{\theta}} = \mathbb{E}_{\pi_{\theta}} \left[\sum_{h=0}^{\infty} \gamma^{h} r_{h} \right] \\ + \eta \nabla_{\theta} V^{\pi_{\theta}} |_{\theta = \theta_{t}} \end{array}$$

Policy Gradient Theorem [Sutton, McAllester, Singh, Mansour]:

$$(s, a) := Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$$
, we have:

Adjust the probability $\pi_{\theta}(a \mid s)$ proportional to $A^{\pi_{\theta}}(s, a) := Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$

Consider tabular MDPs, wit

th
$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}, \ \theta_{s,a} \in \mathbb{R}$$

Consider tabular MDPs, with

 $\frac{\partial V(\theta)}{\partial \theta_{s,a}} = \frac{1}{1-\gamma} d^{\pi}_{\mu}(s) \pi_{\theta}(a \mid s) A^{\pi_{\theta}}(s, a), \text{ where } A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$

th
$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}, \ \theta_{s,a} \in \mathbb{R}$$

PG formulation:

Consider tabular MDPs, with

PG formulation: $\frac{\partial V(\theta)}{\partial \theta_{s,a}} = \frac{1}{1-\gamma} d^{\pi}_{\mu}(s) \pi_{\theta}(a \mid s) A^{\pi_{\theta}}(s, a), \text{ where } A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$

th
$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}, \ \theta_{s,a} \in \mathbb{R}$$

Despite being non-concave, we have global convergence:

Consider tabular MDPs, with

$$\frac{\partial V(\theta)}{\partial \theta_{s,a}} = \frac{1}{1 - \gamma} d^{\pi}_{\mu}(s) \pi_{\theta}(a \mid s) A^{\pi_{\theta}}(s)$$

Despite being non-concave, we have global convergence:

Theorem (Informal) [Agarwal, Kakade, Lee, Mahajan 20; Mei, Xiao, Szepesvari, Schuurmans 20].

Assume $\mu(s) > 0, \forall s$, the PG algorithm $\theta^{t+1} := \theta^t + \eta \nabla_{\theta} V(\theta) |_{\theta = \theta^t}$ converges to global optimality

th
$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}, \ \theta_{s,a} \in \mathbb{R}$$

PG formulation:

s, a), where $A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$

[Kakade 03]

$$F_{\theta} = \mathbb{E}_{s,a \sim d^{\pi_{\theta}}} \left[\nabla_{\theta} \ln \pi_{\theta}(a) \right]$$

[Kakade 03]

Define Fisher information matrix

 $u(s) \left(\nabla_{\theta} \ln \pi_{\theta}(a|s) \right)^{\mathsf{T}} \in \mathbb{R}^{d_{\theta} \times d_{\theta}}$

$$F_{\theta} = \mathbb{E}_{s,a \sim d^{\pi_{\theta}}} \left[\nabla_{\theta} \ln \pi_{\theta}(a) \right]$$

 $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta) |_{\theta = \theta^t}$

[Kakade 03]

Define Fisher information matrix

 $|s) \left(\nabla_{\theta} \ln \pi_{\theta}(a \,|\, s) \right)^{\mathsf{T}} \in \mathbb{R}^{d_{\theta} \times d_{\theta}}$

Natural policy gradient uses F_{θ} to pre-condition PG:

$$F_{\theta} = \mathbb{E}_{s,a \sim d^{\pi_{\theta}}} \left[\nabla_{\theta} \ln \pi_{\theta}(a) \right]$$

 $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta) |_{\theta = \theta^t}$

(For simplicity, assume F_{θ} is full rank —- otherwise use pseudo inverse)

[Kakade 03]

Define Fisher information matrix

 $|s| \left(\nabla_{\theta} \ln \pi_{\theta}(a \,|\, s) \right)^{\mathsf{T}} \in \mathbb{R}^{d_{\theta} \times d_{\theta}}$

Natural policy gradient uses F_{θ} to pre-condition PG:



NPG as a Trust-region optimization procedure: $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta=\theta^{t}} \rangle, \text{ s.t., } KL \left(\rho_{\theta^{t}} | | \rho_{\theta} \right) \leq \delta$ [Bagnell & Schneider 03]

$$\rangle, \text{ s.t., } KL\left(\rho_{\theta^{t}} | | \rho_{\theta}\right) \leq \delta \\ \left(\rho_{\theta}(\tau) := \mu(s_{0}) \prod_{h} \pi(a_{h} | s_{h}) P(s_{h+1} | h_{h}) \right)$$



NPG as a Trust-region $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta = \theta^{t}} \rangle$

i.e., optimize the **linearized objective** s.t. a KL constraint **forcing new policy's trajectory distribution staying close to old one's**

[Bagnell & Schneider 03]

NPG as a Trust-region optimization procedure:

$$angle, ext{ s.t., } KL\left(
ho_{ heta^t} | \, | \,
ho_{ heta}
ight) \leq \delta$$

$$\left(\rho_{\theta}(\tau) := \mu(s_0) \prod_h \pi(a_h \mid s_h) P(s_{h+1} \mid s_h) \right)$$



 $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta = \theta^{t}} \rangle$

i.e., optimize the linearized objective s.t. a KL constraint forcing new policy's trajectory distribution staying close to old one's

[Bagnell & Schneider 03]

NPG as a Trust-region optimization procedure:

$$angle, ext{ s.t., } KL\left(
ho_{ heta^t} | \, | \,
ho_{ heta}
ight) \leq \delta$$

$$\left(\rho_{\theta}(\tau) := \mu(s_0) \prod_h \pi(a_h \mid s_h) P(s_{h+1} \mid s_h) \right)$$

Further perform second-order Taylor expansion on $KL(\rho_{\theta^t} | | \rho_{\theta})$ at θ^t :



NPG as a Trust-region $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta = \theta^{t}} \rangle$

i.e., optimize the **linearized objective** s.t. a KL constraint **forcing new policy's trajectory distribution staying close to old one's**

Further perform second-order Tay $KL\left(
ho_{ heta^t}|\,|\,
ho_{ heta}
ight)pprox 0$

[Bagnell & Schneider 03]

NPG as a Trust-region optimization procedure:

$$angle, ext{ s.t., } KL\left(
ho_{ heta^t} | \, | \,
ho_{ heta}
ight) \leq \delta$$

$$\left(\rho_{\theta}(\tau) := \mu(s_0) \prod_h \pi(a_h \mid s_h) P(s_{h+1} \mid s_h) \right)$$

Further perform second-order Taylor expansion on $KL(\rho_{\theta^t} | | \rho_{\theta})$ at θ^t :

 $KL\left(\rho_{\theta^{t}} | | \rho_{\theta}\right) \approx (\theta - \theta^{t})^{\mathsf{T}} F_{\theta^{t}}(\theta - \theta^{t})$



NPG as a Trust-region $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta = \theta^{t}} \rangle$

i.e., optimize the **linearized objective** s.t. a KL constraint **forcing new policy's trajectory distribution staying close to old one's**

Further perform second-order Tages $KL\left(
ho_{ heta^t}|\,|\,
ho_{ heta}
ight)pprox KL\left(
ho_{ heta^t}|\,|\,
ho_{ heta}
ight)pprox KL\left(
ho_{ heta^t}|\,|\,
ho_{ heta}
ight)$

NPG then is revealed by solving the convex program:

 $\max_{\theta} \langle \theta, \nabla_{\theta} V(\theta) |_{\theta = \theta^{t}} \rangle, \text{ s.t., } (\theta - \theta^{t})^{\mathsf{T}} F_{\theta^{t}}(\theta - \theta^{t}) \leq \delta$

[Bagnell & Schneider 03]

NPG as a Trust-region optimization procedure:

$$angle, ext{ s.t., } KL\left(
ho_{ heta^t} | \, | \,
ho_{ heta}
ight) \leq \delta$$

$$\left(\rho_{\theta}(\tau) := \mu(s_0) \prod_h \pi(a_h \mid s_h) P(s_{h+1} \mid s_h) \right)$$

- Further perform second-order Taylor expansion on $KL(\rho_{\theta^t} | | \rho_{\theta})$ at θ^t :
 - $KL\left(\rho_{\theta^{t}} | | \rho_{\theta}\right) \approx (\theta \theta^{t})^{\mathsf{T}} F_{\theta^{t}}(\theta \theta^{t})$



Recall the softmax Policy for Tabular MDPs:

 $\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A$

$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

Recall the softmax Policy for Tabular MDPs:

 $\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A$

$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

We can show that the NPG update $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta^t)$ is equivalent to (see the exercise in recitation):

Recall the softmax Policy for Tabular MDPs:

$$\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A \quad \pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

$$(\pi^t := \pi_{\theta^t}) \pi^{t+1}(a \mid s) \propto \pi^t(a \mid s) \cdot \exp\left(\eta A^{\pi^t}(s, a)\right)$$

We can show that the NPG update $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta^t)$ is equivalent to (see the exercise in recitation):

Recall the softmax Policy for Tabular MDPs:

$$\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A \quad \pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

 $(\pi^t := \pi_{\theta^t}) \pi^{t+1}(a \mid s) \propto \pi^t(a \mid s)$

Proof sketch: $A^{\pi_{\theta^t}}(\cdot, \cdot) \propto \arg \min \|\nabla_{\theta} V(\theta^t) - F_{\theta^t} x\|_2^2$ (see recitation for details) X

We can show that the NPG update $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta^t)$ is equivalent to (see the exercise in recitation):

$$a \mid s) \cdot \exp\left(\eta A^{\pi^t}(s,a)\right)$$

Recall the softmax Policy for Tabular MDPs:

$$\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A \quad \pi_{\theta}(a \mid s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

$$(\pi^t := \pi_{\theta^t}) \pi^{t+1}(a \mid s) \propto \pi^t(a \mid s) \cdot \exp\left(\eta A^{\pi^t}(s, a)\right)$$

Proof sketch: $A^{\pi_{\theta^t}}(\cdot, \cdot) \propto \arg \min \|\nabla_{\theta} V(\theta^t) - F_{\theta^t} x\|_2^2$ (see recitation for details) \mathcal{X}

We can show that the NPG update $\theta^{t+1} := \theta^t + \eta F_{\theta^t}^{-1} \nabla_{\theta} V(\theta^t)$ is equivalent to (see the exercise in recitation):

Interpretation: for each state s, NPG runs online mirror ascent with $A^{\pi'}(s, \cdot) \in \mathbb{R}^{|A|}$ as the reward vector at iter t

Global Convergence of the exact Natural policy gradient

(Note here we are studying the idealized case where we have exact $A^{\pi'}(\cdot, \cdot)$). We will look into learning/approximation in the recitation)

 $\pi^{t+1}(a \mid s) \propto \pi^t(a \mid s) \cdot \exp\left(\eta A^{\pi^t}(s, a)\right)$

Global Convergence of the exact Natural policy gradient

 $\pi^{t+1}(a \mid s) \propto \pi^t(a \mid s)$

Theorem [Agarwal, Kakade, Lee, Mahaja iterations, there exits a

 $V^{\pi} > V^{\star} - V$

$$a | s) \cdot \exp\left(\eta A^{\pi^t}(s, a)\right)$$

(Note here we are studying the idealized case where we have exact $A^{\pi'}(\cdot, \cdot)$). We will look into learning/approximation in the recitation)

an 20]: Initialize
$$\pi^0(\cdot | s) = \text{Unif}(A)$$
. After T
policy $\pi \in \{\pi^0, ..., \pi^{T-1}\}$, s.t.,
 $\log A = \frac{1}{(1 - \gamma)^2 T}$.

Global Convergence of the exact Natural policy gradient

 $\pi^{t+1}(a \mid s) \propto \pi^t(a \mid s)$

Theorem [Agarwal, Kakade, Lee, Mahaja iterations, there exits a

 $V^{\pi} > V^{\star} - \overline{}$

$$a | s) \cdot \exp\left(\eta A^{\pi^t}(s, a)\right)$$

(Note here we are studying the idealized case where we have exact $A^{\pi'}(\cdot, \cdot)$). We will look into learning/approximation in the recitation)

an 20]: Initialize
$$\pi^0(\cdot | s) = \text{Unif}(A)$$
. After T
policy $\pi \in \{\pi^0, ..., \pi^{T-1}\}$, s.t.,
 $\frac{\log A}{\eta T} = \frac{1}{(1-\gamma)^2 T}$.

Global optimality despite non-concavity in the objective

• No |S| dependence at all; log-dependence on |A|

• No coverage requirement on the initial distribution μ

1. Since we run Mirror Ascent per state, we have that for all $s \in S$:

1. Since we run Mirror Ascent per state, we have that for all $s \in S$:

$$\sum_{t=0}^{T-1} \langle \pi^{\star}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle - \underbrace{\langle \pi^{t}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle}_{0} \lesssim \sqrt{\ln(|A|)T}.$$

regret of mirror ascent on s

=0

1. Since we run Mirror Ascent per state, we have that for all $s \in S$:

$$\sum_{t=0}^{T-1} \langle \pi^{\star}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle - \underbrace{\langle \pi^{t}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle}_{=0} \leq \sqrt{\ln(|A|)T}.$$

regret of mirror ascent on s

$$\sum_{t=0}^{T-1} V^{\pi^{\star}} - V^{\pi^{t}} \propto \sum_{t=0}^{T-1} \mathbb{E}_{s \sim d^{\pi^{\star}}_{\mu}} \left[\mathbb{E}_{a \sim \pi^{\star}(\cdot|s)} A^{\pi^{t}}(s,a) \right] \lesssim \sqrt{\ln(|A|)T}.$$

=0

2. Add $\mathbb{E}_{s \sim d_u^{\pi^*}}$ on both sides, and via performance difference lemma [Kakade & Langford 2003]:
Proof Sketch for NPG's global optimality (a $1/\sqrt{T}$ rate)

1. Since we run Mirror Ascent per state, we have that for all $s \in S$:

$$\sum_{t=0}^{T-1} \langle \pi^{\star}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle - \underbrace{\langle \pi^{t}(\cdot \mid s), A^{\pi^{t}}(s, \cdot) \rangle}_{=0} \lesssim \sqrt{\ln(|A|)T}$$

regret of mirror ascent on s

$$\sum_{t=0}^{T-1} V^{\pi^{\star}} - V^{\pi^{t}} \propto \sum_{t=0}^{T-1} \mathbb{E}_{s \sim d^{\pi^{\star}}_{\mu}} \left[\mathbb{E}_{a \sim \pi^{\star}(\cdot|s)} A^{\pi^{t}}(s,a) \right] \lesssim \sqrt{\ln(|A|)T}.$$

(see the exercise in recitation for a detailed proof with approximation on $Q^{\pi^{*}}$, and see chapter 10 in AJKS monograph for the proof for 1/T rate)

=0

2. Add $\mathbb{E}_{s \sim d_{u}^{\pi^{\star}}}$ on both sides, and via performance difference lemma [Kakade & Langford 2003]:

Policy Gradient and NPG:

Global Convergence vanilla PG and NPG in tabular MDPs with softmax parameterization

NPG w/ approximation in Recitation

Summary so far:

Part 1C: Exploration in tabular MDP w/ UCB-Value Iteration

In this part:

Question: how to explore efficient if we do not know (P, r)

We need to perform efficient exploration when learning:

The combination lock problem:



We need to perform efficient exploration when learning:

The combination lock problem:



The prob of a random walk reaching the goal is exponentially small wrt ${\cal H}$

We need to perform efficient exploration when learning:

The combination lock problem:



The prob of a random walk reaching the goal is exponentially small wrt ${\cal H}$

The principle behind UCB-VI: Optimism in the face of uncertainty

Setting: episodic finite horizon tabular MDP (horizon = H), fixed initial state s_0

transitions $\{P_h\}_{h=0}^{H-1}$ unknown, but reward r(s, a) known

learning protocol:

Goal:

transitions $\{P_h\}_{h=0}^{H-1}$ unknown, but reward r(s, a) known

learning protocol:

1. Learner initializes a policy π^0

Goal:

Setting: episodic finite horizon tabular MDP (horizon = H), fixed initial state s_0

Setting: episodic finite horizon tabular MDP (horizon = H), fixed initial state s_0

transitions $\{P_h\}_{h=0}^{H-1}$ unknown, but reward r(s, a) known

learning protocol:

1. Learner initializes a policy

2. At episode n, learner executes π^n to draw a trajectory starting at s_0 : $\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$, with $a_h^n = \pi^n(s_h^n), r_h^n = r(s_h^n, a_h^n), s_{h+1}^n \sim P(\cdot \mid s_h^n, a_h^n)$

Goal:

y
$$\pi^0$$

Setting: episodic finite horizon tabular MDP (horizon = H), fixed initial state s_0

transitions $\{P_h\}_{h=0}^{H-1}$ unknown, but reward r(s, a) known

learning protocol:

1. Learner initializes a policy

2. At episode n, learner executes π^n to draw a trajectory starting at s_0 : $\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$, with $a_h^n = \pi^n(s_h^n), r_h^n = r(s_h^n, a_h^n), s_{h+1}^n \sim P(\cdot \mid s_h^n, a_h^n)$

3. Learner updates policy to π^{n+1} using all prior information

Goal:

y
$$\pi^0$$

Setting: episodic finite horizon tabular MDP (horizon = H), fixed initial state s_0

transitions $\{P_h\}_{h=0}^{H-1}$ unknown, but reward r(s, a) known

learning protocol:

1. Learner initializes a policy

2. At episode n, learner executes π^n to draw a trajectory starting at s_0 : $\{s_h^n, a_h^n, r_h^n\}_{h=0}^{H-1}$, with $a_h^n = \pi^n(s_h^n), r_h^n = r(s_h^n, a_h^n), s_{h+1}^n \sim P(\cdot \mid s_h^n, a_h^n)$

3. Learner updates policy to π^{n+1} using all prior information

Goal:

Sub-linear regret:

$$\mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] = \operatorname{poly}(S, A, H)\sqrt{N}$$

y
$$\pi^0$$

Use all previous data to estimate transitions $\widehat{P}_{0}^{n}, \ldots, \widehat{P}_{H-1}^{n}$

- Use all previous data to estimate transitions $\widehat{P}_{0}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$

- Use all previous data to estimate transitions $\widehat{P}_{0}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$
- Optimistic planning with learned model: $\pi^n = \text{Value-Iter}\left(\{\widehat{P}_h^n, r_h + b_h^n\}_{h=1}^{H-1}\right)$

- Use all previous data to estimate transitions $\widehat{P}_{0}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$
- Optimistic planning with learned model: $\pi^n = \text{Value-Iter}\left(\{\widehat{P}_h^n, r_h + b_h^n\}_{h=1}^{H-1}\right)$
- Collect a new trajectory by executing π^n in the real world $\{P_h\}_{h=0}^{H-1}$ starting from s_0

Let us consider the **very beginning** of episode *n*:

$$\mathcal{D}_h^n = \{s_h^i\}$$

 $\{a_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$

Let us consider the **very beginning** of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

Let us consider the **very beginning** of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h, \quad N_h^n(s,a,s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s,a,s')\}, \forall s, h \in \mathbb{N}\}$$



Let us consider the **very beginning** of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h, \quad N_h^n(s,a,s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s,a,s')\}, \forall s, h \in \mathbb{N}\}$$

Estimate model $\widehat{P}_{h}^{n}(s'|s,a), \forall s,a,s',h$ (i.e., MLE):

$$\widehat{P}_{h}^{n}(s'|s,a) = \frac{N_{h}^{n}(s,a,s')}{N_{h}^{n}(s,a)}$$



UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $b_h^n(s,a) = ch$

$$H_{h} = \frac{\ln(SAHN/\delta)}{N_{h}^{n}(s,a)}$$

UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$



UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$



UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $\widehat{V}_{H}^{n}(s) = 0, \forall s$



UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $b_h^n(s,a) = ch$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$

$$H_{1} = \frac{\ln(SAHN/\delta)}{N_{h}^{n}(s,a)}$$

Encourage to explore new state-actions

UCBVI – Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $b_h^n(s,a) = c$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$
$$\widehat{V}_{h}^{n}(s) = \max_{a} \ \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \ \widehat{Q}_{h}^{n}(s,a), \forall s$$

$$H_{1} \frac{\ln(SAHN/\delta)}{N_{h}^{n}(s,a)}$$

Encourage to explore new state-actions

UCBVI—Part 2: Reward Bonus Design and Value Iteration

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $b_h^n(s,a) = c$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$

$$\widehat{V}_{h}^{n}(s) = \max_{a} \ \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \ \widehat{Q}_{h}^{n}(s,a), \forall s \qquad \left\| \begin{array}{c} \widehat{V}_{h}^{n} \\ \end{array} \right\|_{\infty} \leq H, \forall s \in \mathbb{N}$$

$$H_{1} \frac{\ln(SAHN/\delta)}{N_{h}^{n}(s,a)}$$

Encourage to explore new state-actions



UCBVI: Put All Together

For
$$n = 1 \to N$$
:
1. Set $N_h^n(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h$
2. Set $N_h^n(s, a, s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s, a, s')\}, \forall s, a, a', h$
3. Estimate \widehat{P}^n : $\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall s, a, s', h$
4. Plan: $\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right), \text{ with } b_h^n(s, a) = cH\sqrt{\frac{\ln(SAHN/\delta)}{N_h^n(s, a)}}$
5. Execute $\pi^n : \{s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n\}$

Theorem: UCBVI Regret Bound

We will prove the following in the recitation:

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}\sqrt{S^{2}AN}\right)$$

Theorem: UCBVI Regret Bound

We will prove the following in the recitation:

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}\sqrt{S^{2}AN}\right)$$

Note that we consider expected regret here (policy π^n is a random quantity). High probability version is not hard to get (need to do a martingale argument)

Remarks:

Theorem: UCBVI Regret Bound

We will prove the following in the recitation:

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}\sqrt{S^{2}AN}\right)$$

Remarks:

- Note that we consider expected regret here (policy π^n is a random quantity). High probability version is not hard to get (need to do a martingale argument)

Dependency on H and S are suboptimal; but the same algorithm can achieve $H^2\sqrt{SAN}$ in the leading term [Azar et.al 17 ICML]

Key Intuition behind the theorem:

VI at episode n under
$$\{\widehat{P}_{h}^{n}\}_{h}$$
 and $\{r_{h} + b_{h}^{n}\}_{h}$
 $\widehat{V}_{H}^{n}(s) = 0, \forall s \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$
 $\widehat{V}_{h}^{n}(s) = \max_{a} \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \widehat{Q}_{h}^{n}(s,a), \forall s$

Key Intuition behind the theorem:

VI at episode n under
$$\{\widehat{P}_{h}^{n}\}_{h}$$
 and $\{r_{h} + b_{h}^{n}\}_{h}$
 $\widehat{V}_{H}^{n}(s) = 0, \forall s \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot | s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$
 $\widehat{V}_{h}^{n}(s) = \max_{a} \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \widehat{Q}_{h}^{n}(s,a), \forall s$

Key lemma 1: optimism – our bonus is large enough s.t. $\widehat{V}_h^n(s) \ge V_h^{\star}(s), \forall s, h$

Key Intuition behind the theorem:

VI at episode n under
$$\{\widehat{P}_{h}^{n}\}_{h}$$
 and $\{r_{h} + b_{h}^{n}\}_{h}$
 $\widehat{V}_{H}^{n}(s) = 0, \forall s \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot | s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$
 $\widehat{V}_{h}^{n}(s) = \max_{a} \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \widehat{Q}_{h}^{n}(s,a), \forall s$

Key lemma 2: regret decomposition:

Regret at iter $n = V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le 1$

Key lemma 1: optimism — our bonus is large enough s.t. $\widehat{V}_h^n(s) \ge V_h^{\star}(s), \forall s, h$

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0})$$
Key Intuition behind the theorem:

VI at episode n under
$$\{\widehat{P}_{h}^{n}\}_{h}$$
 and $\{r_{h} + b_{h}^{n}\}_{h}$
 $\widehat{V}_{H}^{n}(s) = 0, \forall s \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot | s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$
 $\widehat{V}_{h}^{n}(s) = \max_{a} \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \widehat{Q}_{h}^{n}(s,a), \forall s$

Key lemma 2: regret decomposition:

Regret at iter $n = V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le$ $\leq \sum_{h} \mathbb{E}_{s,a \sim d_h^{\pi^n}}$

Key lemma 1: optimism — our bonus is large enough s.t. $\widehat{V}_h^n(s) \ge V_h^{\star}(s), \forall s, h$

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0})$$

$$\left[b_{h}^{n}(s, a) + (\widehat{P}_{h}^{n}(\cdot | s, a) - P_{h}^{\star}(\cdot | s, a))^{\mathsf{T}} \widehat{V}_{h+1}^{n}\right]$$

Key Intuition behind the theorem:

VI at episode n under
$$\{\widehat{P}_{h}^{n}\}_{h}$$
 and $\{r_{h} + b_{h}^{n}\}_{h}$
 $\widehat{V}_{H}^{n}(s) = 0, \forall s \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot | s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$
 $\widehat{V}_{h}^{n}(s) = \max_{a} \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \widehat{Q}_{h}^{n}(s,a), \forall s$

Key lemma 2: regret decomposition:

Regret at iter $n = V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}$ $\leq \sum_{h} \mathbb{E}_{s,a \sim d_h^{\pi^n}}$

If π^n is suboptimal, i.e., $V^*(s_0) - V^{\pi^n}(s_0)$ is large, then π^n must visit some (s, a)pairs with large bonus b(s, a) or wrong $\widehat{P}(\cdot | s, a)$

Key lemma 1: optimism — our bonus is large enough s.t. $\widehat{V}_{h}^{n}(s) \geq V_{h}^{\star}(s), \forall s, h$

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0})$$

$$\left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot | s, a) - P_{h}^{\star}(\cdot | s, a))^{\mathsf{T}} \widehat{V}_{h+1}^{n} \right]$$

Summary

1. Basics of MDPs:

2. Policy Gradient:

3. Efficient exploration in tabular MDPs:

Bellman Equation / Optimality; two planning algs: Value Iteration and Policy Iteration

Vanilla PG formulation & Natural Policy Gradient with their global convergence

The UCB-VI algorithm via the principle of optimism in the face of uncertainty