

Policy Gradient

Problem Setup: • want to optimize $J(\pi) := \mathbb{E} \left[\sum_{t=1}^H \gamma^{t-1} r_t | \pi \right]$

- Large state space. Parametrized $\pi_\theta \leftarrow$

example: $|A|=2$. given state features $\phi(s)$.

$$\pi_\theta(s) = \begin{cases} -1, & \phi(s)^\top \theta > 0 \\ 1, & \phi(s)^\top \theta \leq 0. \end{cases}$$

deterministic policy.
below we mostly deal
w/ stoch poling.

- $\max_{\theta} J(\pi_\theta)$. (SGD: calculate/estimate $\nabla_{\theta} J(\pi_\theta)$)

- Central Question: How to estimate $\nabla_{\theta} J(\pi_\theta)$.

w/o knowing MDP. but using data?

- YES! We can form an unbiased estimate of $\nabla_{\theta} J(\pi_\theta)|_{\theta=\theta_0}$. given an on-policy trajectory.

Data: $s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H$.

where $s_t, a_t \sim \pi_\theta$.

Goal: derive a func of this traj. whose expectation = $\nabla_{\theta} J(\pi_\theta)$

Derivation. Notations & assumptions.

1. $\mathcal{T} = (s_1, a_1, s_2, a_2, s_3, a_3, \dots, s_H, a_H)$.
2. Assume wlog. deterministic reward.
3. Finite state/action spaces.

$$4. R(\mathcal{T}) = \sum_{t=1}^H \gamma^{t-1} R(s_t, a_t).$$

$$\begin{aligned}
 \nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} \mathbb{E}[R(\tau) \mid \pi_{\theta}] \\
 &= \nabla_{\theta} \left(\sum_{\tau} \underbrace{P^{\pi_{\theta}}(\tau)}_{\text{depends on } \theta} \cdot \underbrace{R(\tau)}_{\text{constant!}} \right) \\
 &= \sum_{\tau} \nabla_{\theta} (P^{\pi_{\theta}}(\tau) \cdot R(\tau)) = \sum_{\tau} R(\tau) \cdot \nabla_{\theta} P^{\pi_{\theta}}(\tau). \\
 &= \sum_{\tau} R(\tau) \cdot P^{\pi_{\theta}}(\tau) \cdot \nabla_{\theta} \log P^{\pi_{\theta}}(\tau). \\
 \boxed{P^{\pi_{\theta}}(\tau) = d_{\theta}(s_1) \cdot \pi_{\theta}(a_1 | s_1) \cdot P(s_2 | s_1, a_1) \cdot \pi_{\theta}(a_2 | s_2) \cdot \dots \cdot \pi_{\theta}(a_H | s_H)} \\
 &= \sum_{\tau} R(\tau) \cdot \cancel{P^{\pi_{\theta}}(\tau)} \cdot \nabla_{\theta} \left(\cancel{\log d_{\theta}(s_1)} + \cancel{\log \pi_{\theta}(a_1 | s_1)} + \cancel{\log P(s_2 | s_1, a_1)} \right. \\
 &\quad \left. + \dots + \cancel{\log \pi_{\theta}(a_H | s_H)} \right). \\
 &= \sum_{\tau} \cancel{P^{\pi_{\theta}}(\tau)} \left(R(\tau) \cdot \sum_{t=1}^H \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \\
 &= \mathbb{E} \left[\quad \mid \pi_{\theta} \right].
 \end{aligned}$$

Naïve PG estimator (REINFORCE)

$$\frac{R(\tau)}{\cancel{\pi}} \cdot \sum_{t=1}^H \nabla_{\theta} \log \pi_{\theta}(a_t | s_t).$$

Remarks: (1) Similarity w/ IS derivations — coincidence?

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta \theta \rightarrow 0} \frac{\overbrace{J(\pi_{\theta+\Delta \theta}) - J(\pi_{\theta})}^{\Delta J}}{\Delta \theta}$$

- w/ traj. from π_θ . can estimate
- $J(\pi_\theta)$: on-policy MC.
 - $J(\pi_{\theta+\Delta\theta})$: per-traj. importance sampling
 - step-wise LS ? \rightarrow

Alternatively: $\nabla_\theta J(\pi_\theta) = \nabla_\theta \sum_{\tau} P^{\pi_\theta}(\tau) \cdot R(\tau)$.

$$\text{improved ver: } \nabla_\theta J(\pi_\theta) = \nabla_\theta \sum_{t=1}^H \gamma^{t-1} \mathbb{E}[r_t | \pi_\theta].$$

$$= \sum_{t=1}^H \gamma^{t-1} \nabla_\theta \mathbb{E}[r_t | \pi_\theta].$$

$$= \sum_{t=1}^H \gamma^{t-1} \nabla_\theta \sum_{\tau_{1:t}} P^{\pi_\theta}(\tau_{1:t}) \cdot r_t$$

= ... (leave as exercise).

$\tau_{1:t}$ is the t-step prefix of τ .
 $s_1, a_1, s_2, a_2, \dots, s_t, a_t$.

"vanilla PG".

$$= \mathbb{E} \left[\sum_{t=1}^H \left(\sum_{t'=t}^H \gamma^{t'-t} r_{t'} \right) \cdot \nabla \log \pi_\theta(a_t | s_t) \mid \pi_\theta \right].$$

Actor-critic "form of PG"

$$\sum_{t=1}^H \gamma^{t-1} \mathbb{E} \left[\left(\sum_{t'=t}^H \gamma^{t'-t} r_{t'} \right) \cdot \nabla \log \pi_\theta(a_t | s_t) \mid \pi_\theta \right]$$

$$= \sum_{t=1}^H \gamma^{t-1} \mathbb{E} \left[\mathbb{E} \left[\dots \mid \dots \mid \pi_\theta, s_t, a_t \right] \right]$$

function of s_t, a_t .

$$\nabla \log \pi_\theta(a_t | s_t) \underbrace{\mathbb{E} \left[\sum_{t'=t}^H \gamma^{t'-t} r_{t'} \right]}_{Q^\pi(s_t, a_t)} \Big| s_t, a_t$$

therefore: $\nabla_\theta J(\pi_\theta) = \sum_{t=1}^H \gamma^{t-1} \underbrace{\mathbb{E}_{s_t, a_t \sim \pi_\theta} [Q^\pi(s_t, a_t) \cdot \nabla \log \pi_\theta(a_t | s_t)]}_{\Delta}$

- REINFORCE / Vanilla PG is on-policy.



Actor-Critic: $\nabla_\theta J(\pi_\theta) \approx \frac{1}{1-\gamma} \mathbb{E}_{s_t, a_t \sim \pi_\theta} [\hat{Q}^\pi(s_t, a_t) \cdot \nabla \log \pi_\theta(a_t | s_t)]$

- \hat{Q}^π can be estimated separately. via e.g. TD.
[and can be off-policy.]
- still require on-policy data to approximate $\mathbb{E}_{s_t, a_t \sim \pi_\theta} [\cdot]$

- In general π_θ has to be stochastic.

Parametrized Stochastic policy.

A popular example: linear + softmax.

- Assume we know $\phi: S \times A \rightarrow \mathbb{R}^d$.
- $\pi_\theta(a|s) \propto e^{\theta^\top \phi(s, a)}$ $\Leftrightarrow \pi_\theta(a|s) = \boxed{\frac{e^{\theta^\top \phi(s, a)}}{\sum_{a'} e^{\theta^\top \phi(s, a')}}}$
- "proportional to"

$$\begin{aligned}
 \nabla_{\theta} \log \pi_{\theta}(a|s) &= \nabla \log \frac{e^{\theta^T \phi(s, a)}}{\sum a' e^{\theta^T \phi(s, a')}} \\
 \Rightarrow \left(\log e^{\theta^T \phi(s, a)} - \log \sum_{a'} (\dots) \right) &= \nabla(\theta^T \phi(s, a)) - \nabla(\log \sum_{a'} (\dots)) \\
 &= \phi(s, a) - \nabla \log \left(\sum_{a'} e^{\theta^T \phi(s, a')} \right) \\
 &= \phi(s, a) - \frac{\sum_{a'} e^{\theta^T \phi(s, a')} \cdot \phi(s, a')}{\sum_{a'} e^{\theta^T \phi(s, a')}} \quad \leftarrow \\
 &\qquad\qquad\qquad \xrightarrow{\mathbb{E}_{a' \sim \pi_{\theta}(\cdot|s)} [\phi(s, a')]}
 \end{aligned}$$

Remark:

- Compare to softmax policy in on-policy control w/ SARSA.
w/ linear func. approx of $\mathbb{Q}^{\pi} \approx \theta^T \phi(s, a)$
softmax policy: $\pi(a|s) \propto e^{\theta^T \phi(s, a)/T}$.
- Why T in SARSA and not here?
→ in SARSA: if $\theta^{*\top} \phi(s, a) = \mathbb{Q}^{\pi}(s, a)$
then can't use $2\theta^*$
→ in PG: can arbitrarily rescale θ .
- When linear softmax policies can represent π^* ?
 $\pi_{\theta}(a|s) \propto e^{\boxed{\phi^T(s, a) \theta}}$ → doesn't have to be a value func.
 to represent π^* { ① Assume $\mathbb{Q}^* = \phi^T(s, a) \cdot \theta^*$.
 ② Take $\underline{\theta = C \cdot \theta^*}$ w/ large C .
 in the limit of $C \rightarrow +\infty \Rightarrow \pi_{\theta} \rightarrow \pi^*$.

- Variance reduction in PG. (w/o introducing bias).

$$\nabla J(\pi_\theta) = \frac{1}{1-\gamma} \mathbb{E}_{\substack{(s,a) \sim d^{\pi_\theta}}} \left[\left(Q^{\pi_\theta}(s,a) - f(s) \right) \nabla \log \pi_\theta(a|s) \right].$$

"state baseline".

Unbiasedness: fixing arbitrary state s .

$$\mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [f(s) \cdot \nabla \log \pi_\theta(a|s)].$$

$$= f(s) \cdot \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\nabla \log \pi_\theta(a|s)].$$

$$= \sum_a \cancel{\pi_\theta(a|s)} \frac{\nabla \pi_\theta(a|s)}{\cancel{\pi_\theta(a|s)}}$$

$$= \nabla \sum_a \pi_\theta(a|s) = \nabla(1) = \vec{0}.$$

A special case: $f(s) := V^\pi(s)$. "advantage".

$$\nabla J(\pi_\theta) = \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi_\theta}} \left[(Q^\pi(s,a) - V^\pi(s)) \nabla \log \pi_\theta(a|s) \right].$$