

# Policy Gradient

VI 4 PI can yield  $\pi^t$  and  $\pi^{t+1}$  that are very different, leading to instabilities  
 $\mathcal{Q}$  &  $\mathcal{V}$  can be close but  $\pi \neq \hat{\pi}$  can be far.

$\Rightarrow$  change policy incrementally by parametrizing  $\pi_\theta$  + doing gradient ascent.

$\pi : S \rightarrow A$  can deal with continuous state spaces.

$\rightarrow \pi_\theta(a|s)$   $\min_{\theta} TE[V_{\pi^*}(s)] - \sum_{t=0}^{\infty} E[V_{\pi_{\theta_t}}(s)] =: J(\pi_\theta)$

$\left[ \begin{array}{l} \theta_0 \text{ initiative} \\ \theta_{t+1} \leftarrow \theta_t + \eta \nabla_{\theta} J(\pi_{\theta_t}) \end{array} \right.$

Issue 1:  $\pi_\theta$  needs to be differentiable ( $\neq$  deterministic policies)  
 stochastic policies:

eg. softmax policies  $\pi_\theta(a|s) \propto \exp(\theta_{s,a})$

loglinear "  $\pi_\theta(a|s) \propto \exp(\theta \cdot \phi_{s,a})$

neural softmax policies  $\pi_\theta(a|s) \propto \exp(f_\theta(s,a))$

Issue 2:  $\nabla_{\theta} J(\pi_{\theta_t})$

$\hookrightarrow$  may require differentiation through entire dynamics - JMDP.

## Policy Gradient Theorem (REINFORCE version)

$$\begin{aligned} \nabla_{\theta} E[V_{\pi_{\theta}}(s)] &= \nabla_{\theta} E_{\tau \sim P_{d_0}^{\pi_{\theta}}} [R(\tau)] & R(\tau) &= \sum_{h=1}^H r_h \\ &= \nabla_{\theta} \sum_{\tau} R(\tau) Pr_{d_0}^{\pi_{\theta}}(\tau) & \tau &= (s_1, a_1, r_1, s_2, \dots, s_H, a_H, r_H) \\ &= \sum_{\tau} R(\tau) \nabla_{\theta} Pr_{d_0}^{\pi_{\theta}}(\tau) & \nabla x &= x \nabla \log x \\ &= \sum_{\tau} R(\tau) Pr_{d_0}^{\pi_{\theta}}(\tau) \nabla \log Pr_{d_0}^{\pi_{\theta}}(\tau) \\ &= \sum_{\tau} R(\tau) Pr_{d_0}^{\pi_{\theta}}(\tau) \nabla \log [d_0(s_1) \pi_{\theta}(a_1|s_1) P(s_2|s_1, a_1) \\ & \quad \dots \pi_{\theta}(a_H|s_H) P(s_H|s_{H-1}, a_{H-1})] \\ &= \sum_{\tau} R(\tau) Pr_{d_0}^{\pi_{\theta}}(\tau) \left( \sum_{h=1}^H \nabla \log \pi_{\theta}(a_h|s_h) \right) \end{aligned}$$

$$= E_{\pi_{\theta}} [R(c) \sum_{k=1}^H \nabla \log \pi_{\theta}(a_k | s_k)] \leftarrow$$

## REINFORCE

Compute stochastic gradient by 1) sampling trajectories using  $\pi_{\theta}$

2) compute  $R(c) \sum_{k=1}^H \nabla \log \pi_{\theta}(a_k | s_k) \leftarrow$

Note: estimate gradient with accuracy independent of size of state space.

Issue: Variance of gradient estimate can be high.

Per step importance sampling to lower variance

Q-version

$$\nabla_{\theta} J(\pi_{\theta}) = E_{(s,a) \sim d^{\pi_{\theta}}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s,a)] \leftarrow$$

1) generate trajectory using  $\pi_{\theta}$

2) pick random time step  $h$ .

3) compute  $\nabla_{\theta} \log \pi_{\theta}(a_h | s_h) \sum_{k=h}^H r_k$

Policy Gradient Theorem (Q-version)

$$\nabla_{\theta} V_{\pi_{\theta}}(s) = \nabla_{\theta} E_{a \sim \pi_{\theta}(s)} [Q_{\pi_{\theta}}(s,a)] = \nabla_{\theta} \sum_a \pi_{\theta}(a|s) Q_{\pi_{\theta}}(s,a)$$

$$= \sum_a \nabla_{\theta} \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s,a) + \sum_a \pi_{\theta}(a|s) \nabla_{\theta} Q_{\pi_{\theta}}(s,a)$$

$$Q_{\pi_{\theta}}(s,a) = R(s,a) + \sum_{s'} P(s'|s,a) \nabla_{\theta} V_{\pi_{\theta}}(s')$$

$$\nabla_{\theta} Q_{\pi_{\theta}}(s,a) = \sum_{s'} P(s'|s,a) \nabla_{\theta} V_{\pi_{\theta}}(s')$$

$$\Rightarrow \nabla_{\theta} V_{\pi_{\theta}}(s) = \sum_a \nabla_{\theta} \pi_{\theta}(a|s) \cdot Q_{\pi_{\theta}}(s,a) + \sum_a \pi_{\theta}(a|s) \sum_{s'} P(s'|s,a) \nabla_{\theta} V_{\pi_{\theta}}(s')$$

Recursive equation in  $\nabla_{\theta} V_{\pi_{\theta}}(s)$ .

$$\nabla_{\theta} \pi_{\theta} = \pi_{\theta} \nabla_{\theta} \log \pi_{\theta}$$

$$\nabla_{\theta} E_{s \sim d_h^{\pi_{\theta}}} [V_{\pi_{\theta}}(s)] = E_{s \sim d_h^{\pi_{\theta}}, a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q_{\pi_{\theta}}(s,a)]$$

$$+ E_{s \sim d_h^{\pi_{\theta}}, a \sim \pi_{\theta}, s' \sim P(\cdot|s,a)} [\nabla_{\theta} V_{\pi_{\theta}}(s')]$$

\*  $s' \sim d_{h+1}^{\pi_{\theta}}$

$$= \dots = \sum_{k=h}^H E_{s \sim d_k^{\pi_{\theta}}, a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q_{\pi_{\theta}}(s,a)]$$

### Variance reduction in policy gradient

$$\begin{aligned} \Rightarrow E_{a \sim \pi_\theta(s)} \left[ \nabla \log \pi_\theta(a|s) \right] &= \sum_a \pi_\theta(a|s) \nabla \log \pi_\theta(a|s) \\ &= \sum_a \nabla \pi_\theta(a|s) \quad \nabla x = x \nabla \log x \\ &= \nabla \sum_a \pi_\theta(a|s) \\ &= \nabla 1 \\ &= 0 \end{aligned}$$

∴ can add any function  $f: S \rightarrow \mathbb{R}$  to policy gradient as long as  $f$  doesn't depend on  $a/\theta$ .

⇒ doesn't change bias of gradient estimate.

$$Q^{\pi_\theta}(s, a) \xrightarrow{\text{vanilla}} Q^{\pi_\theta}(s, a) - f(s)$$

eg.  $f = V^{\pi_\theta}(s)$   
 ↑  
 "critic"

Actor-critic approach:

$$Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) \equiv A^{\pi_\theta}(s, a) \quad \text{best of both.}$$

↑ actor (policy gradient)      ↑ critic (value iteration / policy iteration)

effective update:

$$\nabla_\theta J(\pi_\theta) = E_{(s,a) \sim d^{\pi_\theta}} \left[ \nabla_\theta \log \pi_\theta(a|s) \left( \underbrace{Q^{\pi_\theta}(s, a) - f(s)}_{A^{\pi_\theta}(s, a) \text{ if } f = V^{\pi_\theta}} \right) \right]$$

### Analysis (sketch):

(stochastic) gradient descent guaranteed to find approx stationary point under mild conditions (even for non-convex)

where size of gradient  $\|\nabla J(\pi_\theta)\|$  is small.

Does small  $\|\nabla J(\pi_\theta)\|$  imply optimality of  $\hat{\pi} = \pi_{\hat{\theta}}$ ?

$$\begin{aligned}
J(\pi^*) - J(\hat{\pi}) &= E_{s \sim d^{\pi^*}} [Q^{\hat{\pi}}(s, \pi^*) - V^{\hat{\pi}}(s)] \quad \text{Performance Diff. lemma.} \\
&\leq E_{s \sim d^{\pi^*}} [Q^{\hat{\pi}}(s, \pi_{Q^{\hat{\pi}}}) - V^{\hat{\pi}}(s)] \\
&\leq \left\| \frac{d\pi^*}{d\hat{\pi}} \right\|_{\infty} E_{s \sim d^{\hat{\pi}}} [Q^{\hat{\pi}}(s, \pi_{Q^{\hat{\pi}}}) - V^{\hat{\pi}}(s)] \\
&= \underbrace{\left\| \frac{d\pi^*}{d\hat{\pi}} \right\|_{\infty}}_{1)} E_{s \sim d^{\hat{\pi}}} \left[ \underbrace{\sum_a Q^{\hat{\pi}}(s, a) (\pi_{Q^{\hat{\pi}}}(a|s) - \hat{\pi}(a|s))}_{2)} \right]
\end{aligned}$$

1) need to ensure  $\hat{\pi}$  covers  $\pi^*$

$\Rightarrow$  need to ensure initialization distribution covers  $\pi^*$   
 $S, \pi$  do

exploration  $S, \pi$  d  
 $\Leftarrow$  uniform

2) Does small  $\|\nabla J(\hat{\pi})\|$  imply small  $\|\pi_{Q^{\hat{\pi}}} - \hat{\pi}\|$ ?

possible under smoothness  
 $\Leftarrow$   $\Leftarrow$  tabular, linear.

$$\text{Regret} \sim \frac{d\pi^*}{d\pi_0} \cdot |S| \cdot |A|$$

### Variance Reduction

- Qversion of REINFORCE
- Actor-critic
- TRPO (Trust Region Policy optimization) - second order
- PPO (Proximal policy optimization) - first order.

PPO-clip

$$\frac{\pi_{\theta}(a|s)}{\pi_{\text{orig}}(a|s)} = \text{clip}$$

$$E \left[ \log \frac{\pi_{\theta}(a|s)}{\pi_{\text{orig}}(a|s)} A_{\pi_{\theta}(s,a)} \right]$$