# Policy Gradient (continue)

# Recap: the REINFORCE Algorithm

$$\tau = \{s_0, a_0, s_1, a_1, \ldots, s_{H-1}, a_{H-1}\}$$

$$\rho_\theta(\tau) = \mu(s_0)\pi_\theta(a_0 \mid s_0)P(s_1 \mid s_0, a_0)\pi_\theta(a_1 \mid s_1)\ldots$$

# Recap: the REINFORCE Algorithm

$$\tau = \{s_0, a_0, s_1, a_1, \ldots, s_{H-1}, a_{H-1}\}$$

$$\rho_\theta(\tau) = \mu(s_0)\pi_\theta(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\pi_\theta(a_1 \,|\, s_1)\ldots$$

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[ \underbrace{\sum_{h=0}^{H-1} r(s_h, a_h)}_{R(\tau)} \right]$$

# Recap: the REINFORCE Algorithm

$$\tau = \{s_0, a_0, s_1, a_1, \ldots, s_{H-1}, a_{H-1}\}$$

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[ \underbrace{\sum_{h=0}^{H-1} r(s_h, a_h)}_{R(\tau)} \right]$$

$$\rho_\theta(\tau) = \mu(s_0)\pi_\theta(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\pi_\theta(a_1 \,|\, s_1)\ldots$$

$$\nabla_\theta J(\pi_\theta)\,|_{\theta=\theta_0} := \mathbb{E}_{\tau \sim \rho_{\theta_0}(\tau)} \left[ \left( \sum_{h=0}^{H-1} \nabla_\theta \ln \pi_{\theta_0}(a_h \,|\, s_h) \right) R(\tau) \right]$$

# Recap: the REINFORCE Algorithm

Initialize $\theta$

While True:

# Recap: the REINFORCE Algorithm

Initialize $\theta$

While True:

Generate $n$ i.i.d trajectories $\tau^1, \ldots, \tau^n$ using $\pi_\theta$

# Recap: the REINFORCE Algorithm

Initialize $\theta$

While True:

    Generate $n$ i.i.d trajectories $\tau^1, \ldots, \tau^n$ using $\pi_\theta$

    Compute gradient: $g = \dfrac{1}{n} \displaystyle\sum_{i=1}^{n} \left( \sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h^i \mid s_h^i) \cdot R(\tau^i) \right)$

State -Action
at traj î. at time
$h$

$E(g)$
$= \nabla_\theta J(\theta)$

# Recap: the REINFORCE Algorithm

Initialize $\theta$

While True:

Generate $n$ i.i.d trajectories $\tau^1, \ldots, \tau^n$ using $\pi_\theta$

$s_0, t_0 \quad s_1, a_1$ $\qquad$ $s_{H-1}, a_{H-1}$

$r(s_n, a_n)$ $\qquad$ $R(\tau)$

Compute gradient: $g = \dfrac{1}{n} \sum_{i=1}^{n} \left( \sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h^i \mid s_h^i) \cdot R(\tau^i) \right)$ $\qquad = \sum_{h=0}^{H-1} r_h$

update: $\theta \Leftarrow \theta + \eta g$ (or adaptive methods like Adam)

# REINFORCE can have high uncertainty

# REINFORCE can have high uncertainty

$$\text{Gradient: } g = \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h^i \mid s_h^i) \cdot R(\tau^i) \right)$$

# REINFORCE can have high uncertainty

Gradient: $g = \dfrac{1}{n} \sum\limits_{i=1}^{n} \left( \sum\limits_{h=0}^{H-1} \nabla_\theta \ln \pi_\theta(a_h^i \mid s_h^i) \cdot R(\tau^i) \right)$

$r(s_h, a_h)$

Often require large n to reduce the variance, especially when policy $\pi_\theta$ is quite random

$s_h\, a_h$

$\nabla_\theta \ln \pi_\theta(a_h \mid s_h)$

$s_h\, a_h$

# Today's Question:

How to reduce Variance in Policy Gradient?

**Outline:**

1. A $Q(s, a)$ based Policy Gradient

2. Variance Reduction via A Baseline

# Notations

$$\mathcal{M} = \{P, r, \gamma, \mu, S, A\} \quad \text{where } s_0 \sim \mu$$

*Intial state Dist*

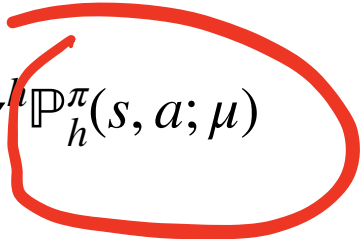$$V^\pi(s) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi\right]$$

Objective: $J(\pi) := \mathbb{E}_{s_0 \sim \mu}\left[V^\pi(s_0)\right]$

# Notations

$$\mathcal{M} = \{P, r, \gamma, \mu, S, A\} \quad \text{where } s_0 \sim \mu$$

$$V^\pi(s) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,|\, s_0 = s, a_h \sim \pi\right]$$

Objective: $J(\pi) := \mathbb{E}_{s_0 \sim \mu}\left[V^\pi(s_0)\right]$

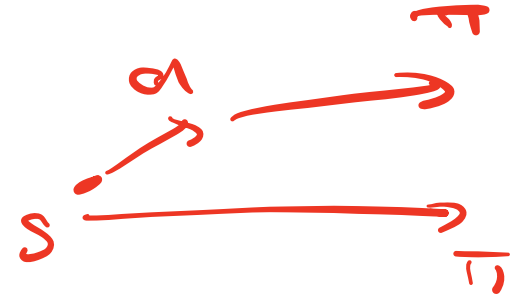$$d^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi(s, a; \mu)$$

# Notations

$$\mathcal{M} = \{P, r, \gamma, \mu, S, A\} \quad \text{where } s_0 \sim \mu$$

$$V^\pi(s) = \mathbb{E}\left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \,|\, s_0 = s, a_h \sim \pi\right]$$

$$\text{Objective: } J(\pi) := \mathbb{E}_{s_0 \sim \mu}\left[V^\pi(s_0)\right]$$

$$d^\pi(s, a) = (1 - \gamma)\sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi(s, a; \mu)$$

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu} \left[ V^{\pi_\theta}(s_0) \right]$$

$$\leftarrow J(\pi_\theta)$$

$$= \mathop{\mathbb{E}}_{s_0 \sim \mu} \left[ \nabla_\theta V^{\pi_\theta}(s_0) \right]$$

# Derivation of Policy Gradient w/ $Q^\pi$

## Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu}\left[V^{\pi_\theta}(s_0)\right]$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\nabla_\theta \underbrace{\mathbb{E}_{a_0 \sim \pi_\theta(s_0)}Q^{\pi_\theta}(s_0, a_0)}_{\triangle}\right]$$

$$V^{\pi_\theta}(s) = \mathop{\mathbb{E}}_{a \sim \pi_\theta(\cdot|s)} Q^{\pi_\theta}(s, a)$$

Chain Rule:

$$\nabla_\theta \left[\left(\sum_a\right) \pi_\theta(a|s_0) Q^{\pi_\theta}(s_0, a_0)\right]$$

$$= \sum_a \left(\nabla_\theta \pi_\theta(a|s_0) Q^{\pi_\theta}(s_0, a) + \pi_\theta(a|s_0) \nabla_\theta Q^{\pi_\theta}(s_0, a)\right)$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu}\left[V^{\pi_\theta}(s_0)\right]$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0)\right] = \mathbb{E}_{s_0 \sim \mu}\left[\sum_{a_0} \nabla_\theta \pi_\theta(a_0 \,|\, s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \,|\, s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0)\right]$$

$\textcircled{2}$

$\textcircled{2}$:

$$\mathbb{E}_{a_0 \sim \pi_\theta(\cdot | s_0)} \nabla_\theta Q^{\pi_\theta}(s_0, a_0)$$

$$= \mathbb{E}_{a_0 \sim \pi_\theta(\cdot | s_0)} \nabla_\theta \left[\Gamma(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P(s_0, a_0)} V^{\pi_\theta}(s_1)\right]$$

$$\nabla_\theta \Gamma(s_0, a_0) = 0$$

$$\gamma: \mathbb{E}_{a_0 \sim \pi_\theta(\cdot | s_0)} \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V(s_1) \quad \boxed{\nabla_\theta V^{\pi_\theta}(s_1)}$$

$\textcircled{1}$

$$\longrightarrow \sum_{a_0} \pi_\theta(a_0 | s_0) \cdot \boxed{\frac{\nabla \pi_\theta(a_0 | s_0)}{\pi_\theta(a_0 | s_0)}} Q^{\pi_\theta}(s_0, a_0)$$

$$= \left(\sum_{a_0} \pi_\theta(a_0 | s_0)\right) \nabla \ln \pi_\theta(a_0 | s_0) \, Q^{\pi_\theta}(s_0, a_0)$$

$$= \mathbb{E}_{a \sim \pi_\theta(\cdot | s_0)} \nabla \ln \pi_\theta(a_0 | s_0) \, Q^{\pi_\theta}(s_0, a_0)$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu} \left[ V^{\pi_\theta}(s_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0) \right] = \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0} \nabla_\theta \pi_\theta(a_0 \mid s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \mid s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0 \in A} \pi_\theta(a_0 \mid s_0) \left[ \frac{\nabla_\theta \pi_\theta(a_0 \mid s_0)}{\pi_\theta(a_0 \mid s_0)} \right] \cdot Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0} \pi_\theta(a_0 \mid s_0) \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V^{\pi_\theta}(s_1) \right]$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu} \left[ V^{\pi_\theta}(s_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0) \right] = \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0} \nabla_\theta \pi_\theta(a_0 \mid s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \mid s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0 \in A} \pi_\theta(a_0 \mid s_0) \left[ \frac{\nabla_\theta \pi_\theta(a_0 \mid s_0)}{\pi_\theta(a_0 \mid s_0)} \right] \cdot Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0} \pi_\theta(a_0 \mid s_0) \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V^{\pi_\theta}(s_1) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \mathbb{E}_{a_0 \sim \pi_\theta(a_0 \mid s_0)} \nabla_\theta \ln \pi_\theta(a_0 \mid s_0) \cdot Q^{\pi_\theta}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_1)$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu}\left[V^{\pi_\theta}(s_0)\right]$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0)\right] = \mathbb{E}_{s_0 \sim \mu}\left[\sum_{a_0} \nabla_\theta \pi_\theta(a_0 \mid s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \mid s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0)\right]$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\sum_{a_0 \in A} \pi_\theta(a_0 \mid s_0)\left[\frac{\nabla_\theta \pi_\theta(a_0 \mid s_0)}{\pi_\theta(a_0 \mid s_0)}\right] \cdot Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0} \pi_\theta(a_0 \mid s_0) \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V^{\pi_\theta}(s_1)\right]$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\mathbb{E}_{a_0 \sim \pi_\theta(a_0 \mid s_0)} \nabla_\theta \ln \pi_\theta(a_0 \mid s_0) \cdot Q^{\pi_\theta}(s_0, a_0)\right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_1)$$

$$= \mathbb{E}_{s_0 \sim \mu}\left[\mathbb{E}_{a_0 \sim \pi_\theta(a_0 \mid s_0)} \nabla_\theta \ln \pi_\theta(a_0 \mid s_0) Q^{\pi_\theta}(s_0, a_0)\right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}}\left[\mathbb{E}_{a_1 \sim \pi_\theta(a_1 \mid s_1)} \nabla_\theta \ln \pi_\theta(a_1 \mid s_1) Q^{\pi_\theta}(s_1, a_1)\right] + \gamma^2 \mathbb{E}_{s_2 \sim \mathbb{P}_2^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_2)$$

← Repeat above

Repeat

$$\gamma^2 \mathbb{E}_{s_2 \sim \mathbb{P}_2^{\pi_\theta}}\left[\mathbb{E}_{a_2 \sim \pi_\theta(\cdot \mid s_2)} \nabla_\theta \ln \pi_\theta(a_2 \mid s_2) Q^{\pi_\theta}(s_2, a_2)\right]$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu} \left[ V^{\pi_\theta}(s_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0) \right] = \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0} \nabla_\theta \pi_\theta(a_0 \,|\, s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \,|\, s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0 \in A} \pi_\theta(a_0 \,|\, s_0) \left[ \frac{\nabla_\theta \pi_\theta(a_0 \,|\, s_0)}{\pi_\theta(a_0 \,|\, s_0)} \right] \cdot Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0} \pi_\theta(a_0 \,|\, s_0) \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V^{\pi_\theta}(s_1) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \mathbb{E}_{a_0 \sim \pi_\theta(a_0 | s_0)} \nabla_\theta \ln \pi_\theta(a_0 \,|\, s_0) \cdot Q^{\pi_\theta}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_1)$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \mathbb{E}_{a_0 \sim \pi_\theta(a_0 | s_0)} \nabla_\theta \ln \pi_\theta(a_0 \,|\, s_0) Q^{\pi_\theta}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \left[ \mathbb{E}_{a_1 \sim \pi_\theta(a_1 | s_1)} \nabla_\theta \ln \pi_\theta(a_1 \,|\, s_1) Q^{\pi_\theta}(s_1, a_1) \right] + \gamma^2 \mathbb{E}_{s_2 \sim \mathbb{P}_2^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_2)$$

$$= \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s_h, a_h \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a_h \,|\, s_h) \cdot Q^{\pi_\theta}(s_h, a_h)$$

$$d^{\pi_\theta} = (1-\gamma) \sum_h \gamma^h \cdot P_h^{\pi_\theta}$$

# Derivation of Policy Gradient w/ $Q^\pi$

Recall definition of value function $V^{\pi_\theta}(s)$

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathbb{E}_{s_0 \sim \mu} \left[ V^{\pi_\theta}(s_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \nabla_\theta \mathbb{E}_{a_0 \sim \pi_\theta(s_0)} Q^{\pi_\theta}(s_0, a_0) \right] = \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0} \nabla_\theta \pi_\theta(a_0 \,|\, s_0) \cdot Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 \,|\, s_0) \cdot \nabla_\theta Q^{\pi_\theta}(s_0, a_0) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \sum_{a_0 \in A} \pi_\theta(a_0 \,|\, s_0) \left[ \frac{\nabla_\theta \pi_\theta(a_0 \,|\, s_0)}{\pi_\theta(a_0 \,|\, s_0)} \right] \cdot Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0} \pi_\theta(a_0 \,|\, s_0) \mathbb{E}_{s_1 \sim P(s_0, a_0)} \nabla_\theta V^{\pi_\theta}(s_1) \right]$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \mathbb{E}_{a_0 \sim \pi_\theta(a_0 | s_0)} \nabla_\theta \ln \pi_\theta(a_0 \,|\, s_0) \cdot Q^{\pi_\theta}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_1)$$

$$= \mathbb{E}_{s_0 \sim \mu} \left[ \mathbb{E}_{a_0 \sim \pi_\theta(a_0 | s_0)} \nabla_\theta \ln \pi_\theta(a_0 \,|\, s_0) Q^{\pi_\theta}(s_0, a_0) \right] + \gamma \mathbb{E}_{s_1 \sim \mathbb{P}_1^{\pi_\theta}} \left[ \mathbb{E}_{a_1 \sim \pi_\theta(a_1 | s_1)} \nabla_\theta \ln \pi_\theta(a_1 \,|\, s_1) Q^{\pi_\theta}(s_1, a_1) \right] + \gamma^2 \mathbb{E}_{s_2 \sim \mathbb{P}_2^{\pi_\theta}} \nabla_\theta V^{\pi_\theta}(s_2)$$

$$= \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s_h, a_h \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a_h \,|\, s_h) \cdot Q^{\pi_\theta}(s_h, a_h) \quad = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a \,|\, s) \cdot Q^{\pi_\theta}(s, a)$$

# Summary so far:

chain rule +  Important weighting + Recursion:

# Summary so far:

chain rule + Important weighting + Recursion:

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s,a \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a_h \mid s_h) \cdot Q^{\pi_\theta}(s_h, a_h)$$

$$R(\tau)$$

# Summary so far:

chain rule +  Important weighting + Recursion:

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s,a \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a_h \,|\, s_h) \cdot Q^{\pi_\theta}(s_h, a_h)$$

for finite horizon MDP (try this out!)

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim \mathbb{P}_h^{\pi_\theta}} \left[ \nabla \ln \pi_\theta(a_h \,|\, s_h) \cdot Q_h^{\pi_\theta}(s_h, a_h) \right]$$
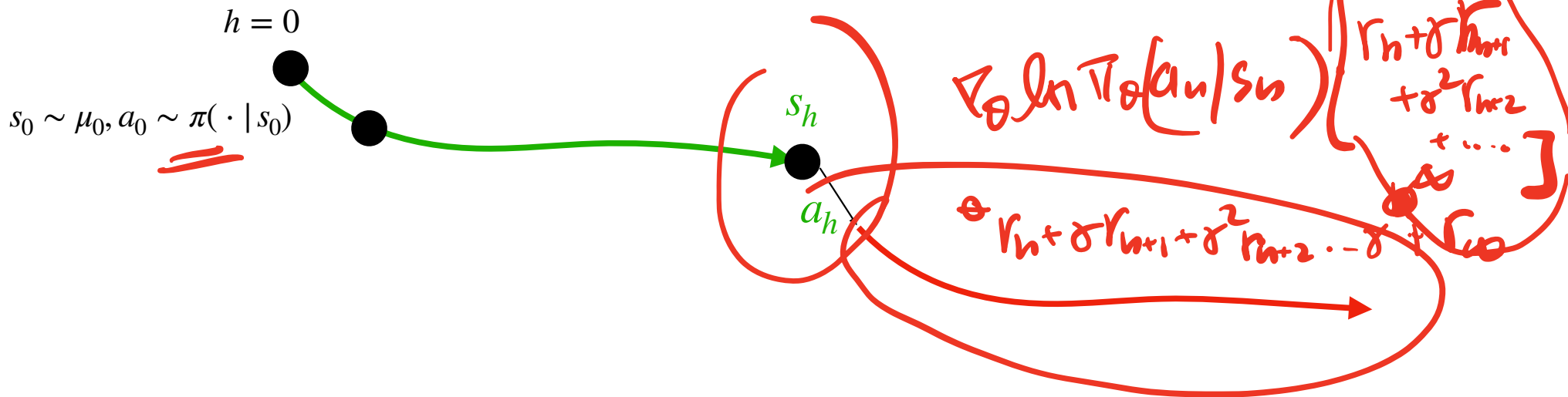
# Summary so far:

chain rule + Important weighting + Recursion:

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s,a \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a_h \mid s_h) \cdot Q^{\pi_\theta}(s_h, a_h)$$

for finite horizon MDP (try this out!)

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{H-1} \mathbb{E}_{s_h, a_h \sim \mathbb{P}_h^{\pi_\theta}} \left[ \nabla \ln \pi_\theta(a_h \mid s_h) \cdot Q_h^{\pi_\theta}(s_h, a_h) \right]$$

$h = 0$

$s_0 \sim \mu_0, a_0 \sim \pi(\cdot \mid s_0)$

$s_h$

$a_h$

$\nabla_\theta \ln \pi_\theta(a_h \mid s_h)$

$\left[ r_h + \gamma r_{h+1} + \gamma^2 r_{h+2} + \dots \right]$

$r_h + \gamma r_{h+1} + \gamma^2 r_{h+2} \dots - \gamma$

**Outline:**

✔ 1. A $Q(s, a)$ based Policy Gradient

2. Variance Reduction via A Baseline

# Q-based PG with a baseline

# Q-based PG with a baseline

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s,a \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a \mid s) \cdot Q^{\pi_\theta}(s, a)$$

# Q-based PG with a baseline

$$\nabla_\theta J(\pi_\theta) = \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s,a \sim \mathbb{P}_h^{\pi_\theta}} \nabla_\theta \ln \pi_\theta(a \mid s) \cdot Q^{\pi_\theta}(s, a)$$

s,a

This still contains the entire future, could still have high variance

# The Advantage-based PG:

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent **baseline, turns out that**

$$b: S \mapsto R$$

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent **baseline, turns out that**

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

$$\nabla_\theta \ln \pi_\theta(a|s) \, b(s) = 0$$

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent ***baseline, turns out that***

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

$= 0, \forall s$

$\mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} \nabla_\theta \ln \pi_\theta(a \mid s) b(s)$

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent **baseline, turns out that**

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

$$\mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} \nabla_\theta \ln \pi_\theta(a \mid s) b(s) \qquad \frac{\nabla \pi_\theta}{\pi_\theta}$$

$$= \sum_a \pi_\theta(a \mid s) \frac{\nabla \pi_\theta(a \mid s)}{\pi_\theta(a \mid s)} b(s)$$

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent **baseline, turns out that**

$$\nabla_\theta J(\pi_\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s,a\sim d_\mu^{\pi_\theta}} \left[ \nabla_\theta \ln \pi_\theta(a \,|\, s) \cdot \left( Q^{\pi_\theta}(s,a) - b(s) \right) \right]$$

$$\mathbb{E}_{a\sim\pi_\theta(\cdot|s)} \nabla_\theta \ln \pi_\theta(a \,|\, s) b(s)$$

$$= \sum_a \pi_\theta(a \,|\, s) \frac{\nabla \pi_\theta(a \,|\, s)}{\pi_\theta(a \,|\, s)} b(s) = b(s) \sum_a \nabla \pi_\theta(a \,|\, s) = b(s) \nabla \left[ \sum_a \pi_\theta(a \,|\, s) \right]$$

$$= ??$$

$$= 1$$

$$\nabla_\theta 1 = 0$$

# The Advantage-based PG:

Denote $b(s)$ as a state-dependent **baseline, turns out that**

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

$$\mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} \nabla_\theta \ln \pi_\theta(a \mid s) b(s)$$

$$= \sum_a \pi_\theta(a \mid s) \frac{\nabla \pi_\theta(a \mid s)}{\pi_\theta(a \mid s)} b(s) \quad = b(s) \sum_a \nabla \pi_\theta(a \mid s) = b(s) \nabla \left[ \sum_a \pi_\theta(a \mid s) \right] = b(s) \nabla 1 = 0$$

# Baseline can further reduce the variance

$$\nabla_\theta J(\pi_\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d_\mu^{\pi_\theta}} \left[ \nabla_\theta \ln \pi_\theta(a \,|\, s) \cdot \left( Q^{\pi_\theta}(s, a) - b(s) \right) \right]$$

# Baseline can further reduce the variance

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

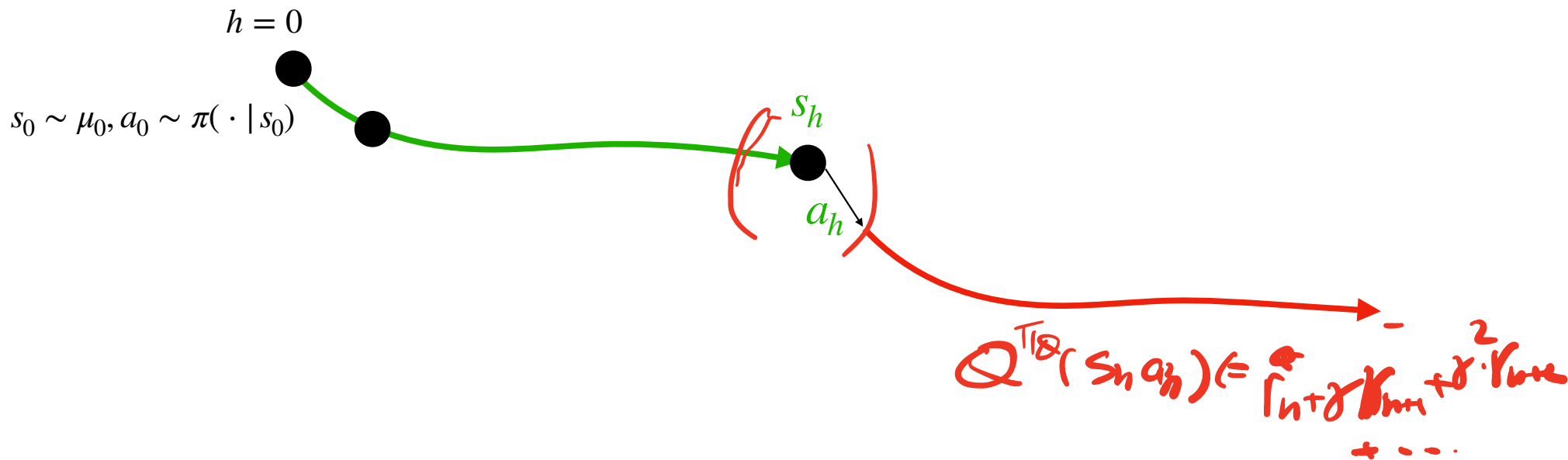Baseline helps reduce variance (formal proof out of scope), and a common choice is $V^{\pi_\theta}(s)$:

$$b(s) = V^{\pi_\theta}(s)$$

# Baseline can further reduce the variance

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

$V^{\pi_\theta}(s)$

Baseline helps reduce variance (formal proof out of scope), and a common choice is $V^{\pi_\theta}(s)$:
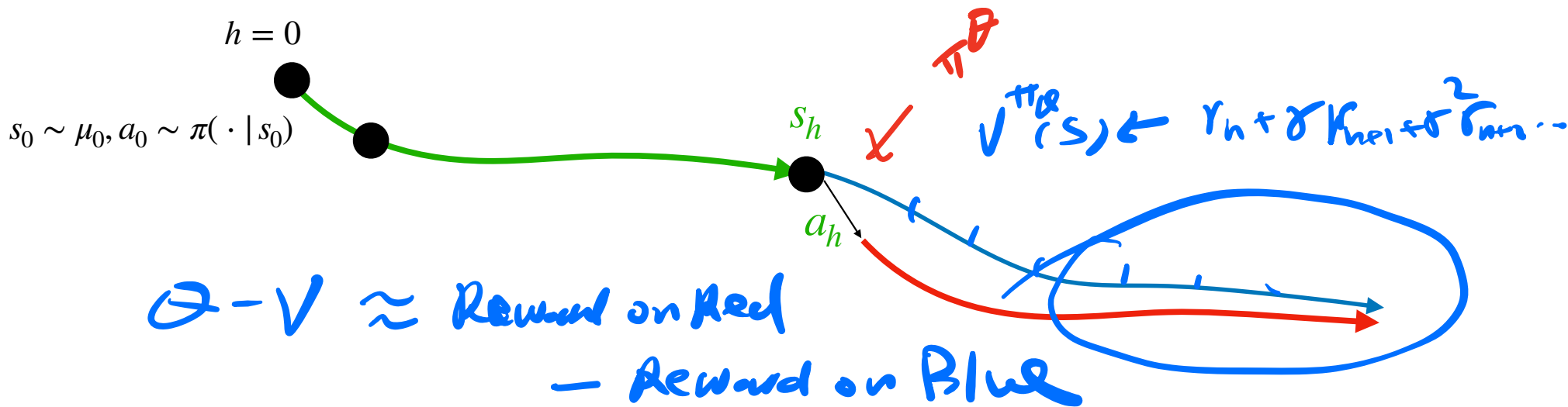
$h = 0$

$s_0 \sim \mu_0, a_0 \sim \pi(\cdot \mid s_0)$

$s_h$

$a_h$

$Q^{\pi\theta}(s_h, a_h) (= \hat{r}_h + \gamma \overline{r_{h+1}} + \gamma^2 \overline{r_{h+2}} + \cdots$

$Q^{\pi}(sa)$

$-V^{\pi}(s) \not\equiv \theta(sa)$

# Baseline can further reduce the variance

$A^{\pi_\theta}(sa)$

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

Baseline helps reduce variance (formal proof out of scope), and a common choice is $V^{\pi_\theta}(s)$:

$h = 0$

$s_0 \sim \mu_0, a_0 \sim \pi(\cdot \mid s_0)$

$s_h$

$\pi^\theta$

$V^{\pi_\theta}(s) \leftarrow r_h + \gamma K_{h+1} + \gamma^2 r_{h+2} \cdots$

$a_h$

$Q - V \approx$ Reward on Red

— Reward on Blue

# Summary for PG:

## Three common PG formulations:

REINFORCE

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[ \left( \sum_{h=0}^{\infty} \nabla_\theta \ln \pi_\theta(a_h \mid s_h) \right) R(\tau) \right]$$

# Summary for PG:

## Three common PG formulations:

### REINFORCE

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[ \left( \sum_{h=0}^{\infty} \nabla_\theta \ln \pi_\theta(a_h \,|\, s_h) \right) R(\tau) \right]$$

### PG w/ $Q$ function

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

# Summary for PG:

## Three common PG formulations:

### REINFORCE

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \rho_\theta(\tau)} \left[ \left( \sum_{h=0}^{\infty} \nabla_\theta \ln \pi_\theta(a_h \mid s_h) \right) R(\tau) \right]$$

### PG w/ $Q$ function

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

### PG w/ $A$ function (use $V^\pi(s)$ as a baseline)

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



$b(s)$
$= V^{\pi_\theta}(s)$