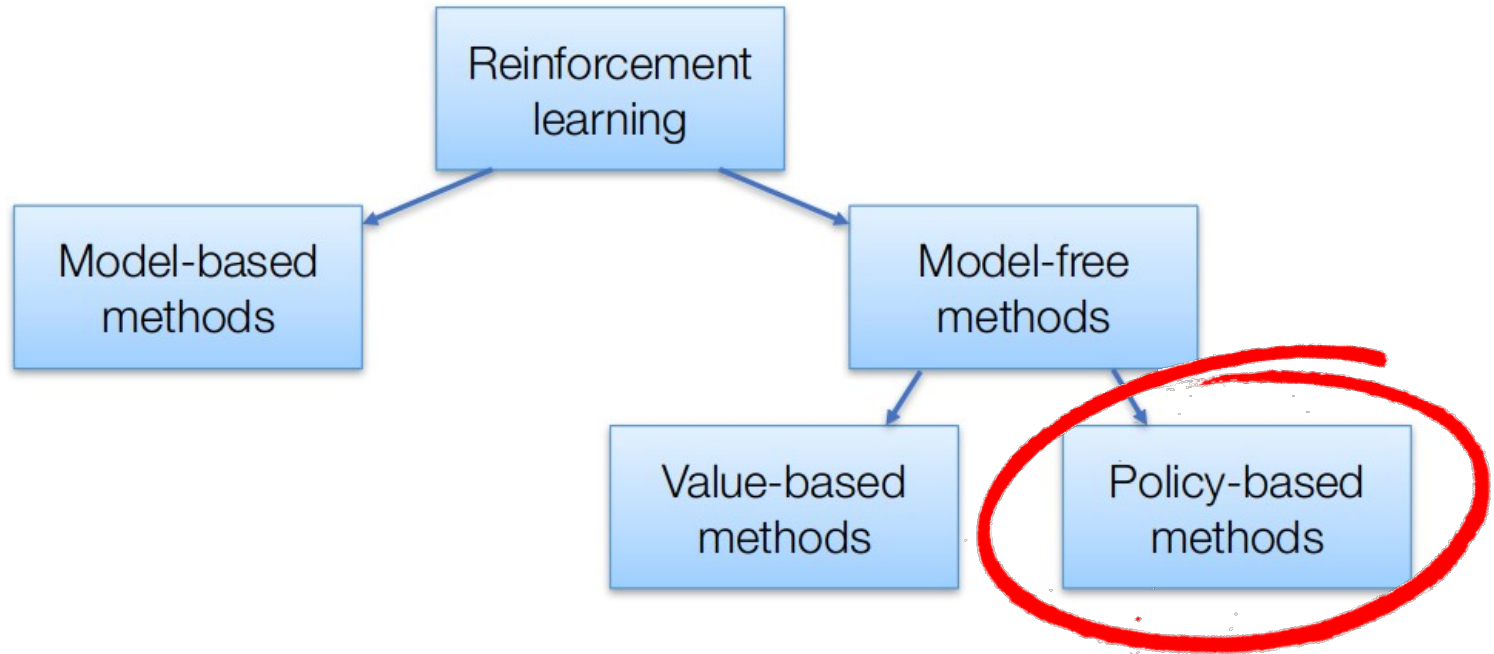


# Policy Gradients and Actor Critic

**Goal:** Find an optimal policy that maximises expected return in a given environment

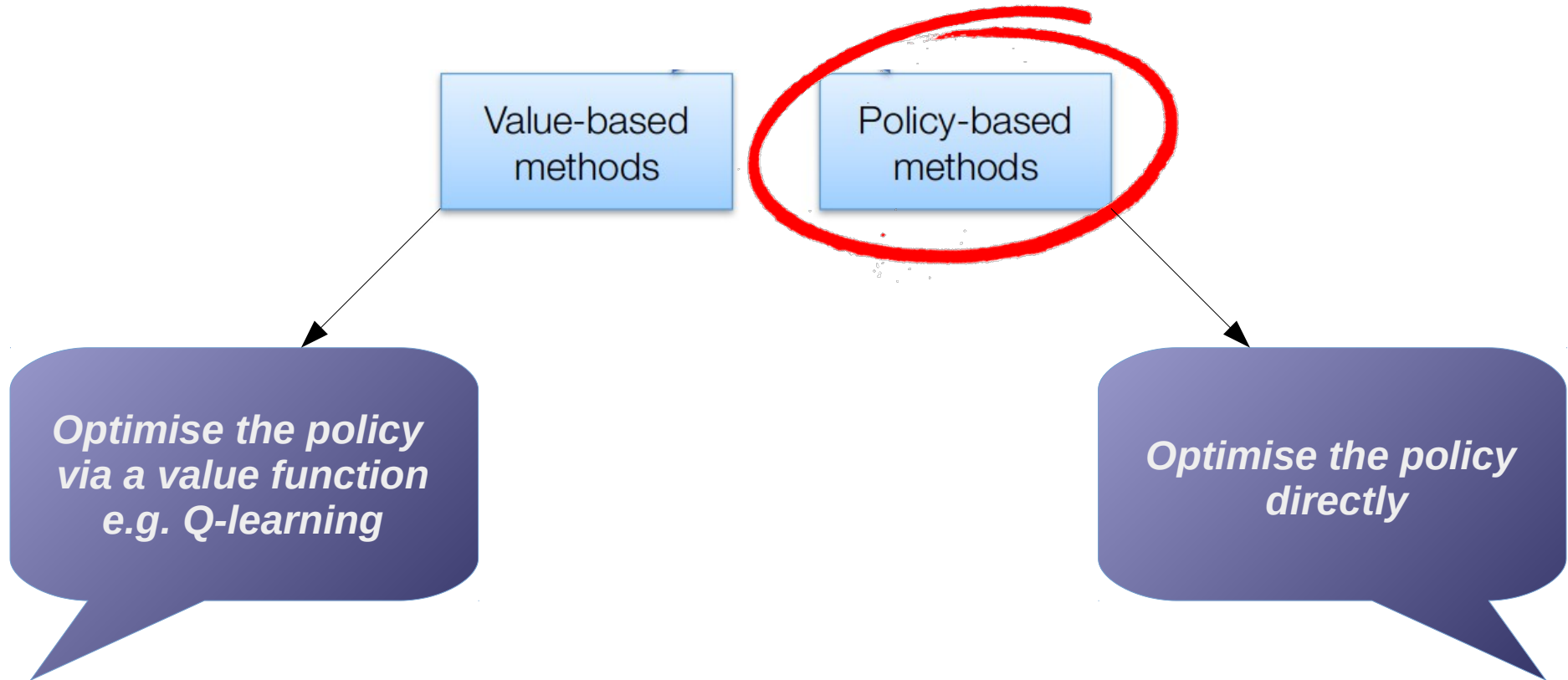
# Policy Gradients

**Goal:** Find an optimal policy that maximises expected return in a given environment



# Policy Gradients








**Goal:** Find an optimal policy that maximises expected return in a given environment



# Policy Gradients

**Goal:** Find an optimal policy that maximises expected return in a given environment

Pros & cons of policy gradient methods (vs Q-learning)

- |   |   |
|---|---|
|  Learns policy directly (good when value function is very complex) |  Can converge to local optima                            |
|  Can learn stochastic policies                                     |  High variance of rewards leads to low sample efficiency |
|  Works with continuous action spaces                               |  Less stable performance                                 |
|  Faster convergence  |   |

# Policy Gradients

**Goal:** Find an optimal policy that maximises expected return in a given environment

## Definitions

$\tau = (s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}, s_H)$  is a trajectory

$\pi_\theta : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is a stochastic policy with parameters  $\theta$

$R(s_t, a_t)$  is the expected reward for taking action  $a_t$  in state  $s_t$

$\mathcal{R}(\tau) = \sum_{t=0}^{T-1} R(s_t, a_t)$  is the expected reward over a trajectory  $\tau$

# Policy Gradients

**Goal:** Find an optimal policy that maximises expected return in a given environment

Now define the Utility function  $U(\theta)$  to be the expected return following policy  $\pi_\theta$

$$\begin{aligned} U(\theta) &= \mathbb{E}_{\pi_\theta}(\mathcal{R}(\tau)) \\ &= \sum_{\tau} P(\tau; \theta) \mathcal{R}(\tau) \end{aligned}$$

**Goal:** Find policy parameters that maximise  $U(\theta)$ :

$$\arg \max_{\theta} (U(\theta)) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) \mathcal{R}(\tau)$$

# Policy Gradients

**Goal:** Find policy parameters that maximise  $U(\theta)$

Gradient ascent:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} [U(\theta)]$

# Policy Gradients

The **REINFORCE** algorithm

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} [U(\theta)]$$

with

$$\begin{aligned} \nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t | s_t))] \right) \mathcal{R}(\tau) \right) \\ &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t^{(i)} | s_t^{(i)})] \right) \mathcal{R}(\tau^{(i)}) \right) \end{aligned}$$



# Policy Gradients

Improvements to the **REINFORCE** algorithm

Recall:  $\theta \leftarrow \theta + \alpha \nabla_{\theta} [U(\theta)]$

Problems:

- High variance of  $\nabla_{\theta} [U(\theta)]$
- (Potentially) slow to converge

# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\nabla_{\theta} [U(\theta)] \approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right)$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{t-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) + \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \cancel{\sum_{t'=1}^{t-1} R(s_{t'}^{(i)}, a_{t'}^{(i)})} + \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right)\end{aligned}$$

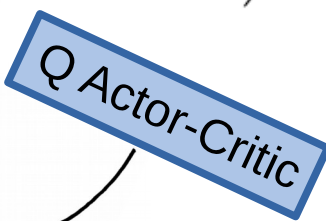
# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \cancel{\sum_{t'=1}^{t-1} R(s_{t'}^{(i)}, a_{t'}^{(i)})} + \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( Q(s_t^{(i)}, a_t^{(i)}) \right) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \mathcal{R}(\tau^{(i)}) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=1}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \cancel{\sum_{t'=1}^{t-1} R(s_{t'}^{(i)}, a_{t'}^{(i)})} + \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( \sum_{t'=t}^{T-1} R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right) \right) \\&= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( Q(s_t^{(i)}, a_t^{(i)}) \right) \right)\end{aligned}$$


# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ : 
$$\nabla_{\theta} [U(\theta)] = \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log(P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b)$$



# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\ &= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\ &= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\ &= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [P(\tau; \theta)]}{P(\tau; \theta)} b \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))]) (\mathcal{R}(\tau) - b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))]) \mathcal{R}(\tau) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))]) b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))]) \mathcal{R}(\tau) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [P(\tau; \theta)]}{P(\tau; \theta)} \right) b \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))]) \mathcal{R}(\tau) \right) - \sum_{\tau} \left( \nabla_{\theta} [P(\tau; \theta)] \right) b\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [P(\tau; \theta)]}{P(\tau; \theta)} b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( \nabla_{\theta} [P(\tau; \theta)] b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} \left[ \sum_{\tau} P(\tau; \theta) \right]\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [P(\tau; \theta)]}{P(\tau; \theta)} b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( \nabla_{\theta} [P(\tau; \theta)] b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} \left[ \sum_{\tau} P(\tau; \theta) \right] \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} [1]\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [P(\tau; \theta)]}{P(\tau; \theta)} b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( \nabla_{\theta} [P(\tau; \theta)] b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} \left[ \sum_{\tau} P(\tau; \theta) \right] \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} [1] \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot 0\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b$ :

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] (\mathcal{R}(\tau) - b)) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] b) \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} P(\tau; \theta) \left( \frac{\nabla_{\theta} [(P(\tau; \theta))]}{P(\tau; \theta)} b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - \sum_{\tau} \left( \nabla_{\theta} [(P(\tau; \theta))] b \right) \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} \left[ \sum_{\tau} \left( (P(\tau; \theta)) \right) \right] \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot \nabla_{\theta} [1] \\&= \sum_{\tau} \left( P(\tau; \theta) (\nabla_{\theta} [\log (P(\tau; \theta))] \mathcal{R}(\tau)) \right) - b \cdot 0 \\&= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) \mathcal{R}(\tau)\end{aligned}$$



# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b(s)$ :

$$\nabla_{\theta} \left[ U(\theta) \right] = \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b(s))$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b(s)$ : e.g  $b(s) = V(s)$

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b(s)) \\ &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - V(s))\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b(s)$ : e.g  $b(s) = V(s)$

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b(s)) \\ &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - V(s)) \\ &= \dots \\ &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} \left[ \log \left( \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \right] \right) \left( Q(s_t^{(i)}, a_t^{(i)}) - V(s_t^{(i)}) \right) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b(s)$ : e.g  $b(s) = V(s)$

$$A(s_t, a_t) = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b(s)) \\ &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - V(s)) \\ &= \dots \\ &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t^{(i)} | s_t^{(i)})] \right) \right) \left( Q(s_t^{(i)}, a_t^{(i)}) - V(s_t^{(i)}) \right) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t^{(i)} | s_t^{(i)})] \right) \right) \left( A(s_t^{(i)}, a_t^{(i)}) \right) \right)\end{aligned}$$

# Policy Gradients

Improvements to the **REINFORCE** algorithm

Introduce a baseline  $b(s)$ : e.g  $b(s) = V(s)$

$$A(s_t, a_t) = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$\begin{aligned}\nabla_{\theta} [U(\theta)] &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - b(s)) \\ &= \mathbb{E}_{\pi_{\theta}} \left( \nabla_{\theta} [\log (P(\tau; \theta))] \right) (\mathcal{R}(\tau) - V(s)) \\ &= \dots\end{aligned}$$

Advantage  
Actor-Critic

$$\begin{aligned}&\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t^{(i)} | s_t^{(i)})] \right) \right) \left( Q(s_t^{(i)}, a_t^{(i)}) - V(s_t^{(i)}) \right) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^{T-1} \left( \nabla_{\theta} [\log (\pi_{\theta}(a_t^{(i)} | s_t^{(i)})] \right) \right) \left( A(s_t^{(i)}, a_t^{(i)}) \right) \end{aligned}$$