

Policy Gradient with Baseline

Shusen Wang

Policy Gradient

- Use policy network, $\pi(a|s; \theta)$, for controlling the agent.
- State-value function:

$$\begin{aligned}\underline{V_\pi(s)} &= \mathbb{E}_{A \sim \pi}[Q_\pi(s, A)] \\ &= \sum_a \underline{\pi(a|s; \theta)} \cdot Q_\pi(s, a).\end{aligned}$$

- Policy gradient:

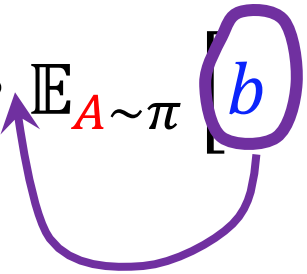
$$\frac{\partial V_\pi(s)}{\partial \theta} = \mathbb{E}_{A \sim \pi} \left[\underline{\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta}} \cdot \underline{Q_\pi(s, A)} \right].$$

Baseline

Baseline

- Let the baseline, b , be anything independent of A .

- $\mathbb{E}_{A \sim \pi} \left[b \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right]$



Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a | s; \theta) \cdot \frac{\partial \ln \pi(a | s; \theta)}{\partial \theta} \end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a | s; \theta) \frac{\partial \ln \pi(a | s; \theta)}{\partial \theta} \end{aligned}$$

$$= \frac{1}{\pi(a | s; \theta)} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a | s; \theta) \cdot \left[\frac{1}{\pi(a | s; \theta)} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \underbrace{\pi(a | s; \theta)}_{\text{cancel}} \cdot \left[\underbrace{\frac{1}{\pi(a | s; \theta)}}_{\text{cancel}} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \cancel{\pi(a | s; \theta)} \cdot \left[\frac{1}{\cancel{\pi(a | s; \theta)}} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \boxed{\sum_a} \frac{\partial \pi(a | s; \theta)}{\partial \theta} \end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned} \bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \cancel{\pi(a | s; \theta)} \cdot \left[\frac{1}{\cancel{\pi(a | s; \theta)}} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \frac{\partial \pi(a | s; \theta)}{\partial \theta} \\ &= b \cdot \frac{\partial \sum_a \pi(a | s; \theta)}{\partial \theta} \end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned}\bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a | s; \theta) \cdot \left[\frac{1}{\pi(a | s; \theta)} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \frac{\partial \pi(a | s; \theta)}{\partial \theta} \\ &= b \cdot \frac{\partial \sum_a \pi(a | s; \theta)}{\partial \theta} \\ &= b \cdot \frac{\partial 1}{\partial \theta}\end{aligned}$$

Baseline

- Let the baseline, b , be anything independent of A .

$$\begin{aligned}\bullet \mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] &= b \cdot \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \pi(a | s; \theta) \cdot \left[\frac{1}{\pi(a | s; \theta)} \cdot \frac{\partial \pi(a | s; \theta)}{\partial \theta} \right] \\ &= b \cdot \sum_a \frac{\partial \pi(a | s; \theta)}{\partial \theta} \\ &= b \cdot \frac{\partial \sum_a \pi(a | s; \theta)}{\partial \theta} \\ &= b \cdot \frac{\partial 1}{\partial \theta} = 0.\end{aligned}$$

Policy Gradient with Baseline

If b is independent of A , then $\mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] = 0$.

- Policy gradient:

$$\begin{aligned} & \frac{\partial V_{\pi}(s)}{\partial \theta} \\ &= \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] \end{aligned}$$

Policy Gradient with Baseline

If b is independent of A , then $\mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] = 0$.

- Policy gradient:

$$\begin{aligned} & \frac{\partial V_{\pi}(s)}{\partial \theta} \\ &= \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] - \underbrace{\mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot b \right]}_{\text{Equal to zero}} \end{aligned}$$

Policy Gradient with Baseline

If b is independent of A , then $\mathbb{E}_{A \sim \pi} \left[b \cdot \frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \right] = 0$.

- Policy gradient:

$$\begin{aligned} & \frac{\partial V_{\pi}(s)}{\partial \theta} \\ &= \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot Q_{\pi}(s, A) \right] - \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot b \right] \\ &= \mathbb{E}_{A \sim \pi} \left[\frac{\partial \ln \pi(A | s; \theta)}{\partial \theta} \cdot (Q_{\pi}(s, A) - b) \right]. \end{aligned}$$

Policy Gradient with Baseline

Theorem. If b is independent of A_t , then policy gradient is equal to:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - b) \right].$$

Monte Carlo Approximation

Monte Carlo Approximation

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - b) \right].$$

- Randomly sample an action $a_t \sim \pi(\cdot | s_t; \theta)$.

Monte Carlo Approximation

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - b) \right].$$

- Randomly sample an action $a_t \sim \pi(\cdot | s_t; \theta)$.
- Compute: $\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, a_t) - b)$.
- $\mathbf{g}(a_t)$ is an unbiased estimate of the policy gradient:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} [\mathbf{g}(A_t)].$$

Monte Carlo Approximation

Stochastic policy gradient with baseline:

$$\mathbf{g}(a_t) = \frac{\partial \ln \pi(a_t | s_t; \theta)}{\partial \theta} \cdot (Q_\pi(s_t, a_t) - b)$$

- Whatever b (independent of A_t) we use, the policy gradient $\mathbb{E}_{A_t \sim \pi}[\mathbf{g}(A_t)]$ remains the same.
- However, b affects the stochastic policy gradient $\mathbf{g}(a_t)$.
- A good b leads to smaller variance and speeds up convergence.

Choices of Baselines

Choice 1: $b=0$

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - b) \right].$$

- We can simply set $b = 0$.
- It becomes the standard policy gradient:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot Q_{\pi}(s_t, A_t) \right].$$

Choice 2: b is state-value

Policy gradient with baseline:

$$\frac{\partial V_{\pi}(s_t)}{\partial \theta} = \mathbb{E}_{A_t \sim \pi} \left[\frac{\partial \ln \pi(A_t | s_t; \theta)}{\partial \theta} \cdot (Q_{\pi}(s_t, A_t) - b) \right].$$

- Because s_t has been observed, $b = V_{\pi}(s_t)$ is independent of A_t .
- Why using such a baseline?
- $V_{\pi}(s_t)$ is close to $Q_{\pi}(s_t, A_t)$:

$$V_{\pi}(s_t) = \mathbb{E}_{A_t} [Q_{\pi}(s_t, A_t)].$$

Thank you!