

Multi-Agent Systems

Introduction to Reinforcement Learning

Part 4: Policy-Gradient and Actor-Critic

Eric Pauwels (CWI & VU)

December 12, 2023

Outline

Policy Gradient Methods

Policy gradient algorithms

Policy gradient algorithms

- optimise policy directly,
- **NOT** via value function (indirectly)

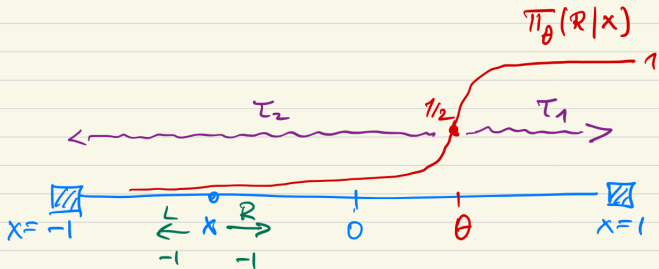
Ingredients:

1. Parametrised policy $\pi_\theta(a|s)$ (θ to be determined)
2. Objective function $J(\theta)$ to be maximised
3. Update rule: $\theta_{new} \leftarrow \theta_{old}$, specifically **gradient ascent**:

$$\theta_{new} \leftarrow \theta_{old} + \nabla_{\theta} J(\theta_{old})$$

Policy gradient: Example(1)

Absorbing states at $x = \pm 1$, absorption reward = 0



$$R(\tau_1) = -(1-\theta)$$

$$p(\tau_1|\theta) = 1/2$$

$$R(\tau_2) = -(1+\theta)$$

$$p(\tau_2|\theta) = 1/2$$

$$R(\tau_2) < R(\tau_1)$$

Policy gradient theorem

- **Objective function** (as path integral and MC version)

$$J(\theta) := \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)] = \int R(\tau) p(\tau | \theta) d\tau \approx \frac{1}{N} \sum_{\tau \sim \pi_\theta} R(\tau)$$

(Path integral makes dependence on θ explicit)

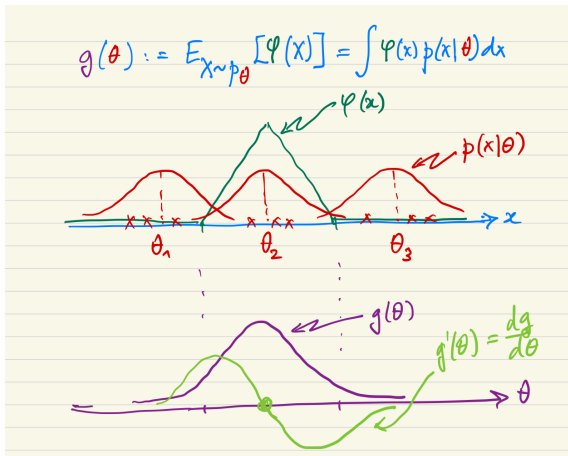
- In **abstract terms** (to simplify notation):

$$g(\theta) := \mathbb{E}_{x \sim p_\theta} [\phi(x)] = \int \phi(x) p(x | \theta) dx$$

- Need to **compute derivative** (to optimise):

$$\frac{d}{d\theta} g(\theta) = \frac{d}{d\theta} \int \phi(x) p(x | \theta) dx = \int \phi(x) \frac{d}{d\theta} (p(x | \theta)) dx$$

Policy gradient theorem



Optimal value $\theta^* = \theta_2$

Policy gradient theorem

- Maximise $g(\theta) = \mathbb{E}_{X \sim p_\theta} [\phi(X)] \approx \frac{1}{n} \sum_{X_i \sim p_\theta} \phi(X_i)$,
- **Intuition:** Sample $X_i \sim p(x | \theta)$
 - If $\phi(X_i)$ **large**: change θ to make X_i **more likely**;
 - If $\phi(X_i)$ **small**: change θ to make X_i **less likely**;
- **Mathematics:** Policy gradient theorem

$$\frac{dg}{d\theta} = \frac{d}{d\theta} \mathbb{E}_{X \sim p_\theta} [\phi(X)] = \mathbb{E}_{X \sim p_\theta} \left[\phi(X) \frac{d}{d\theta} (\log p(X | \theta)) \right]$$

Change θ such that

- If $\phi(X_i)$ **large**: make X_i **more likely**, i.e. $\frac{d(\log)p(X_i)}{d\theta} > 0$
- If $\phi(X_i)$ **small**: make X_i **less likely**, i.e. $\frac{d(\log)p(X_i)}{d\theta} < 0$

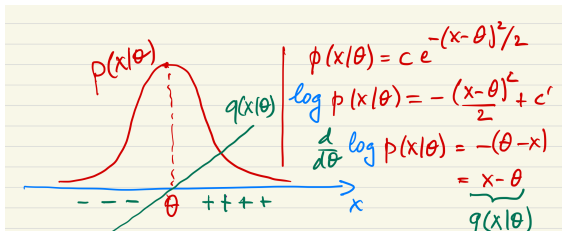
Policy gradient theorem

$$\begin{aligned}
 \frac{d}{d\theta} g(\theta) &= \int \phi(x) \frac{d}{d\theta} (p(x|\theta)) dx \\
 &= \int \phi(x) \left[\frac{\frac{d}{d\theta} (p(x|\theta))}{p(x|\theta)} \right] p(x|\theta) dx \\
 &= \int \phi(x) \left[\frac{\frac{d}{d\theta} (p(x|\theta))}{p(x|\theta)} \right] p(x|\theta) dx \\
 &= \int \phi(x) \left[\frac{d}{d\theta} (\log p(x|\theta)) \right] p(x|\theta) dx \\
 &= \mathbb{E}_{X \sim p_\theta} \left[\phi(X) \frac{d}{d\theta} (\log p(x|\theta)) \right]
 \end{aligned}$$

$$\boxed{\frac{d}{d\theta} \mathbb{E}_{X \sim p_\theta} [\phi(X)] = \mathbb{E}_{X \sim p_\theta} \left[\phi(X) \frac{d}{d\theta} (\log p(X|\theta)) \right]}$$

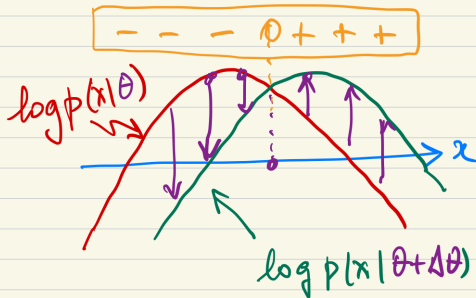
Policy gradient theorem

$$\begin{aligned}
 \frac{d}{d\theta} g(\theta) &= \frac{d}{d\theta} \mathbb{E}_{X \sim p_\theta} [\phi(X)] \\
 &= \mathbb{E}_{X \sim p_\theta} \left[\phi(X) \frac{d}{d\theta} (\log p(X | \theta)) \right] \\
 &\approx \frac{1}{n} \sum_{X_i \sim p_\theta} \phi(X_i) \frac{d}{d\theta} (\log p(X_i | \theta)) \\
 &= \frac{1}{n} \sum_{X_i \sim p_\theta} \phi(X_i) q(X_i | \theta)
 \end{aligned}$$



Policy gradient theorem

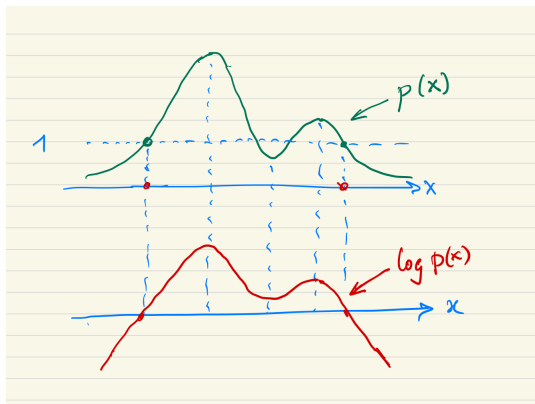
$$\nabla_{\theta} g(\theta) = \mathbb{E}_{x \sim \pi_{\theta}} [\varphi(x) \nabla_{\theta} \log p(x|\theta)]$$



at given x
 how does $\log p(x|\theta)$
 change if θ increases
 i.e. $\theta \rightarrow \theta + \Delta\theta$

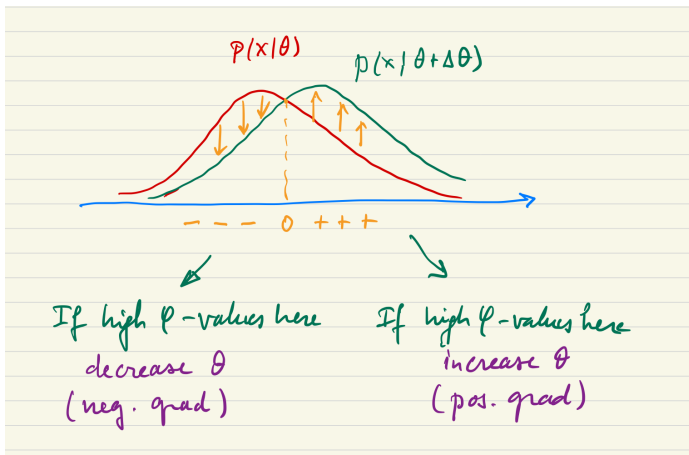
$p(x)$ vs $\log p(x)$

- $p(x)$ vs $\log p(x)$ have same local extremes, increasing and decreasing behaviour
- $\text{sgn}(\nabla_{\theta} p(x | \theta)) = \text{sgn}(\nabla_{\theta} \log p(x | \theta))$

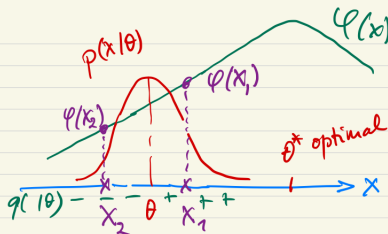


Policy gradient theorem

$$\frac{dg}{d\theta} \approx \frac{1}{n} \sum_{X_i \sim p_\theta} \phi(X_i) \frac{d}{d\theta} (\log p(X_i | \theta))$$

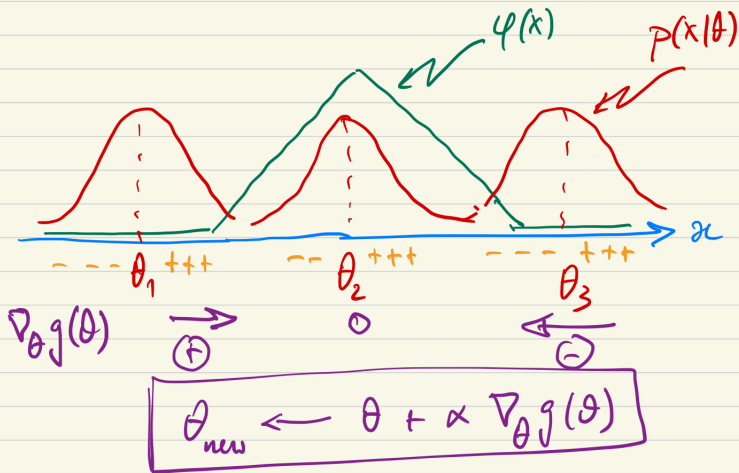


Policy gradient theorem



$$\begin{aligned}
 \nabla_{\theta} q(\theta) &= \mathbb{E}_{x \sim p_{\theta}} [\varphi(x) \nabla_{\theta} \log p(x|\theta)] \\
 &\approx \frac{1}{n} \sum_{x_i \sim p_{\theta}} \varphi(x_i) \underbrace{\nabla_{\theta} \log p(x_i|\theta)}_{q(x_i|\theta)} \\
 &\sim \underbrace{\varphi(x_1) q(x_1|\theta)}_{+} + \underbrace{\varphi(x_2) q(x_2|\theta)}_{-} > 0
 \end{aligned}$$

Policy gradient theorem



Policy gradient theorem

make positive → MAXIMISATION

$$\nabla_{\theta} g(\theta) \approx \frac{1}{n} \sum_{x_i \sim p_{\theta}} \underbrace{\varphi(x_i)}_{\text{change } \theta \text{ s.t.}} \underbrace{\nabla_{\theta} \log p(x_i | \theta)}_{\substack{\text{change } \theta \text{ s.t.} \\ p(x_i | \theta) \text{ increases} \\ (\nabla_{\theta} \log p > 0) \text{ or } \\ p(x_i | \theta) \text{ decreases} \\ (\nabla_{\theta} \log p < 0)}}$$

if $\varphi(x_i)$ HIGH \rightarrow $p(x_i | \theta)$ increases $(\nabla_{\theta} \log p > 0)$

if $\varphi(x_i)$ LOW \rightarrow $p(x_i | \theta)$ decreases $(\nabla_{\theta} \log p < 0)$

Policy gradient theorem

General result

$$\nabla_{\theta} \mathbb{E}_{X \sim p_{\theta}} [\phi(X)] = \mathbb{E}_{X \sim p_{\theta}} [\phi(X) \nabla_{\theta} \log p(X | \theta)]$$

Policy Gradient Theorem

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log p(\tau | \theta)]$$

- If $R(\tau)$ high, change θ to make τ MORE likely, i.e. $p(\tau | \theta) \uparrow$
- If $R(\tau)$ low, change θ to make τ LESS likely, i.e. $p(\tau | \theta) \downarrow$

Policy gradient theorem

$$p(\tau | \theta) = \prod_{t \geq 0} \underbrace{p(s_{t+1} | s_t, a_t)}_{\text{MDP env.}} \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}}$$

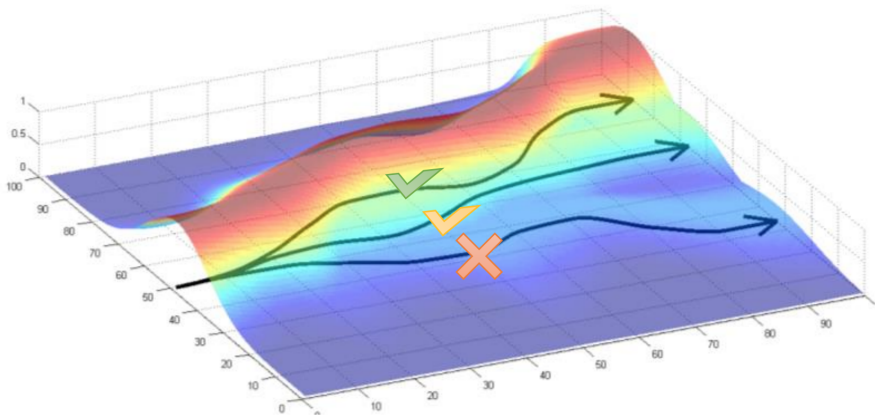
$$\log p(\tau | \theta) = \sum_{t \geq 0} [\log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t)]$$

$$\nabla_{\theta} \log p(\tau | \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

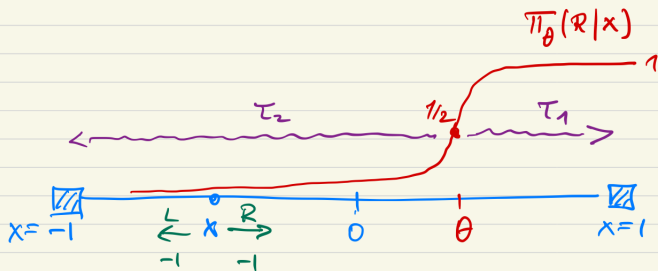
Policy Gradient Theorem

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log p(\tau | \theta)] \\ &= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] \end{aligned}$$

Policy gradient illustration



Policy gradient: Example(1)



$$R(\tau_1) = -(1-\theta)$$

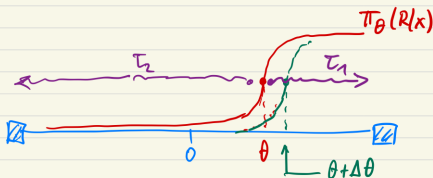
$$p(\tau_1, \theta) = 1/2$$

$$R(\tau_2) = -(1+\theta)$$

$$p(\tau_2|\theta) = 1/2$$

$$R(\tau_2) < R(\tau_1)$$

Policy gradient: Example(2)



$$p(\tau_1|\theta+\Delta\theta) < p(\tau_1|\theta) \rightarrow \nabla_{\theta} \log p(\tau_1|\theta) < 0 \quad (= -\delta)$$

$$p(\tau_2|\theta+\Delta\theta) > p(\tau_2|\theta) \rightarrow \nabla_{\theta} \log p(\tau_2|\theta) > 0 \quad (= +\delta)$$

$$\begin{aligned} \nabla_{\theta} J &\approx R(\tau_1) \cdot \underbrace{\nabla_{\theta} \log p(\tau_1|\theta)}_{-\delta} + R(\tau_2) \underbrace{\nabla_{\theta} \log p(\tau_2|\theta)}_{+\delta} \\ &\approx \delta [-R(\tau_1) + R(\tau_2)] \\ &\approx \delta [|R(\tau_2)| - |R(\tau_1)|] < 0 \Rightarrow \theta \downarrow \end{aligned}$$

Recall: $\text{reward}(\text{path}) = -\text{length}(\text{path})$

REINFORCE algorithm

Example of policy gradient algo

1. Initialise parameter θ of policy π_θ , learning rate α
2. **for** episode = 1 ... NR_EPISODES:
3. Sample trajectory $\tau = \{s_0, a_0, r_1, s_1, \dots, r_T, s_T\}$
4. Set $\nabla_\theta J(\theta) = 0$
5. # add gradient contributions along trajectory
6. **for** $t = 0, 1, \dots, T$:
7. $R_t(\tau) = \sum_{u=t}^T \gamma^{u-t} r_u$,
8. $\nabla_\theta J(\theta) = \nabla_\theta J(\theta) + R_t(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$
9. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$ # update policy parameter

NOTE: REINFORCE algo is **on-policy**:

τ cannot be re-used when policy (i.e. θ) changes!

Worked example

- Linear state space $(-10:10)$, left and right absorbing
- $R_{left} = 0, R_{right} = 5, R_{NT} = -1$

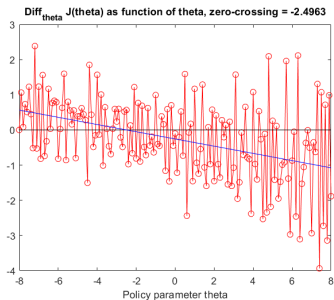
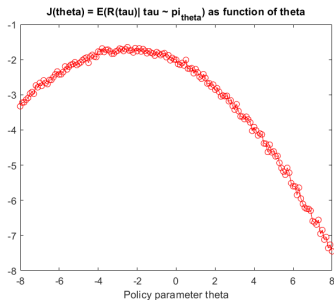


Figure: Left: MC computation for each value of θ , (Right) Numerical gradient (diff), zero crossing at $\theta = -2.5$

Worked example

- Linear state space $(-10:10)$, left and right absorbing
- $R_{left} = 0, R_{right} = 5, R_{NT} = -1$

$$\nabla_{\theta} J(\theta) = \sum_t R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

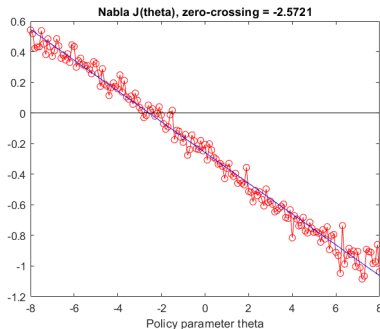


Figure: $\nabla_{\theta} J(\theta)$ via policy gradient theorem, zero crossing at $\theta = -2.57$

Improving REINFORCE algorithm

- Policy gradient estimate has **high variance**
(trajectories might be very different!)

$$\nabla_{\theta} J(\theta) \approx \sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

- **Variance reduction** by introducing **action-independent baseline**:

$$\nabla_{\theta} J(\theta) \approx \sum_{t=0}^T (R_t(\tau) - b(s_t)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

- Example: **Actor-Critic algorithm**:

$$R_t(\tau) = q_{\pi_{\theta}}(s_t, a_t) \quad \text{and} \quad b(s_t) = v_{\pi_{\theta}}(s_t)$$