

ST449 Artificial Intelligence and Deep Learning

Lecture 10

Policy Gradient Methods



Milan Vojnovic

<https://github.com/lse-st449/lectures>

Topics of this lecture

- Policy-based learning
- Policy learning objective functions
- Policy gradient theorem
- Actor-critic algorithms
- Compatible function approximation theorem
- Variance reduction using a baseline

Policies that we studied so far

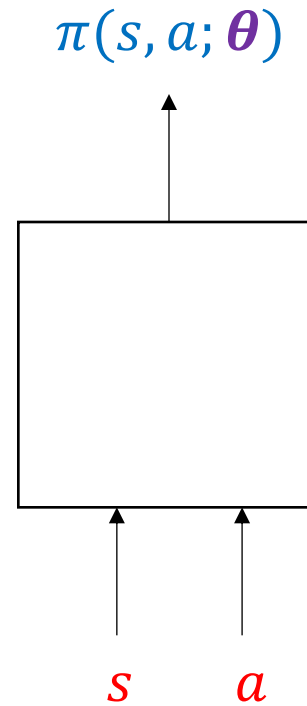
- Greedy:

$$\pi(s, a) = \frac{1}{|\arg \max_{a'} Q(s, a')|} 1_{a \in \arg \max_{a'} Q(s, a')}$$

- Soft policies, e.g. ϵ -greedy:

$$\pi(s, a) = (1 - \epsilon) \frac{1}{|\arg \max_{a'} Q(s, a')|} 1_{a \in \arg \max_{a'} Q(s, a')} + \epsilon \frac{1}{|A(s)|}$$

Policy function approximation



Types of RL with function approximation

Policy function approximation			
		No	Yes
Value function approximation	No		Policy-based
	Yes	Value-based	Actor-critic

Value-based:

- Implicit policy, e.g. ϵ -greedy

Policy-based:

- No value function
- Learn policy

Actor-critic:

- Learn value
- Learn policy

Pros and cons of policy-based learning

Pros:

- Better convergence properties
- Scalable for high-dimensional or continuous action spaces
- Suitable for learning stochastic policies

Cons:

- Convergence to local minima
- Policy evaluation typically inefficient and high variance

Example

A	B	C	D	E
x		£		x

x

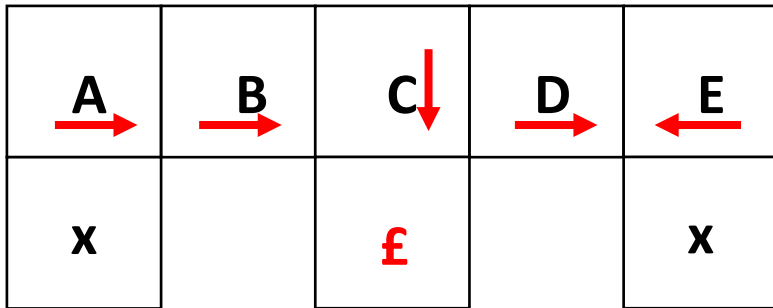
 terminal states

- Suppose we have features vectors:

$$\phi(s, a) = \begin{cases} 1 & \text{if } s \in \{A, B, C, D, E\} \text{ and } a = \text{move to E} \\ 0 & \text{otherwise} \end{cases}$$

- Two approaches:
 - Value-based using the action-value approximator: $\hat{Q}(s, a; \mathbf{w}) = f(\phi(s, a), \mathbf{w})$
 - Policy-based using the policy approximator: $\pi(s, a; \boldsymbol{\theta}) = g(\phi(s, a), \boldsymbol{\theta})$

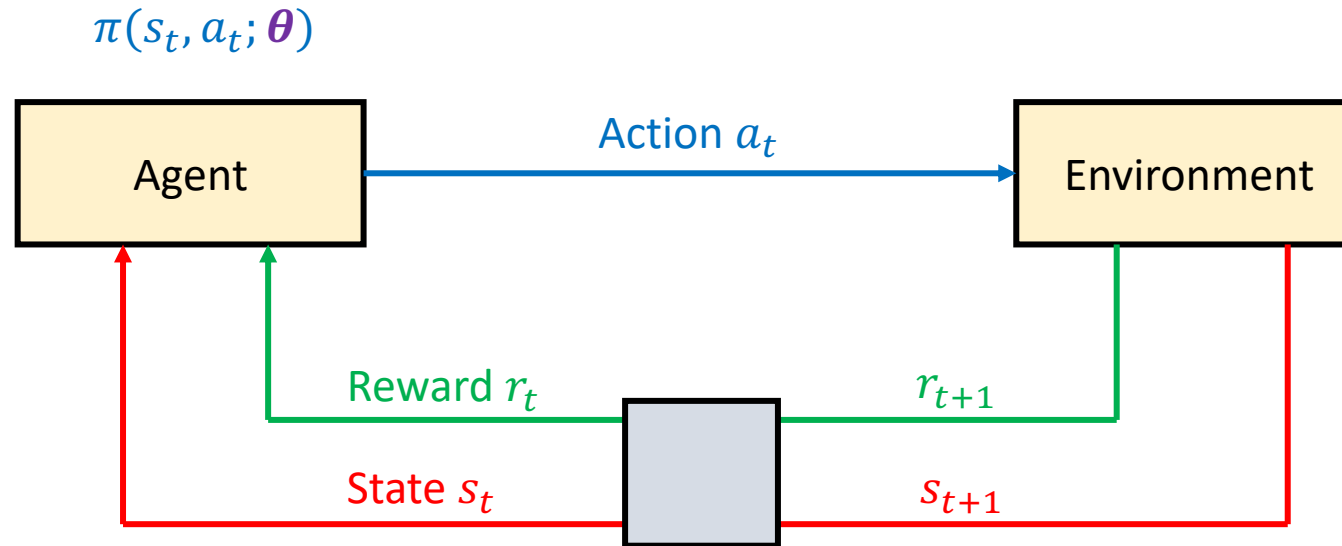
Example (cont'd)



- States **B** and **D** cannot be distinguished because of the symmetry
- A deterministic policy **can get stuck**

- **Optimal stochastic policy**: at states **B** and **D** move either to **E** or **W** equiprobably
- **Value-based learning**: learns a nearly-deterministic policy
 - Long time until reaching **£**
- **Policy-based learning**: can learn an optimal policy

Policy-based reinforcement learning problem



- **Problem:** how to fit parameters of a policy approximator?
- **Solution:** find parameter θ that maximizes an objective function

Policy objective functions

- Average rewards:

$$J(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{E}_{\pi(\cdot, \cdot; \boldsymbol{\theta})} [\sum_{t=1}^n r_t] = \sum_{s,a} \mu^{\pi(\cdot, \cdot; \boldsymbol{\theta})}(s) \pi(s, a; \boldsymbol{\theta}) R_s^a$$

- $R_s^a := \mathbf{E}[r_{t+1} | s_t = s, a_t = a]$
- $\mu^{\pi(\cdot, \cdot; \boldsymbol{\theta})}(\cdot)$ assumed to be the stationary distribution under policy $\pi(\cdot, \cdot; \boldsymbol{\theta})$:

$$\mu^{\pi(\cdot, \cdot; \boldsymbol{\theta})}(s) = \sum_{s', a'} P_{s', s}^{a'} \mu^{\pi(\cdot, \cdot; \boldsymbol{\theta})}(s') \pi(s', a'; \boldsymbol{\theta}) \quad (\text{global balance equations})$$

Policy objective functions (cont'd)

- **Discounted expected rewards:** given a discount factor $\gamma \in [0,1]$ and initial state distribution μ , maximize the expected discounted rewards:

$$J(\theta) = \mathbf{E}_{\pi(\cdot;\theta)}[\sum_{t \geq 0} \gamma^t r_{t+1}]$$

or, equivalently,

$$J(\theta) = \sum_{s^0} \mu(s^0) V^{\pi(\cdot;\theta)}(s^0)$$

- For example, distribution μ may have all its mass on a specific initial state
- If $\gamma = 1$, the task is assumed to be episodic

Policy gradient theorem

- **Thm.** For any differentiable policy $\pi(s, a; \theta)$ with respect to parameter θ the policy gradient for average reward and discounted expected rewards objective is

$$\nabla_{\theta} J(\theta) = \sum_{s,a} \mu^{\pi(\cdot; \theta)}(s, a) \nabla_{\theta} \log(\pi(s, a; \theta)) Q^{\pi(\cdot; \theta)}(s, a)$$

- For average reward objective:

$\mu^{\pi(\cdot; \theta)}$ is the stationary distribution of $\{(s_t, a_t)\}_t$ under policy $\pi(\cdot; \theta)$

- For discounted expected rewards objective:

$$\mu^{\pi(\cdot; \theta)}(s, a) = \sum_{t \geq 0} \gamma^t \sum_{s^0} \Pr_{\pi(\cdot; \theta)} [s_t = s, a_t = a | s_0 = s^0] \mu(s^0)$$

Policy scores

- The term

$$\nabla_{\theta} \log(\pi(s, a; \theta))$$

is often referred as the **policy score**

Comment on discounted expected rewards

- For every $\gamma \in [0,1)$, we have

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{1-\gamma} \mathbf{E}_{(s,a) \sim \tilde{\mu}^{\pi(\cdot;\boldsymbol{\theta})}} \left[\nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta})) Q^{\pi(\cdot;\boldsymbol{\theta})}(s, a) \right]$$

$$\text{where } \tilde{\mu}^{\pi(\cdot;\boldsymbol{\theta})}(s, a; \gamma) = \frac{\mu^{\pi(\cdot;\boldsymbol{\theta})}(s, a; \gamma)}{\sum_{s', a'} \mu^{\pi(\cdot;\boldsymbol{\theta})}(s', a'; \gamma)} = (1-\gamma) \mu^{\pi(\cdot;\boldsymbol{\theta})}(s, a; \gamma)$$

- Abbreviated notation: $\mathbf{E}_{\pi(\cdot;\boldsymbol{\theta})}[\cdot] \equiv \mathbf{E}_{(s,a) \sim \tilde{\mu}^{\pi(\cdot;\boldsymbol{\theta})}}[\cdot]$

Proof for average rewards

- The action-value functions can be defined as:

$$Q^\pi(s, a) = \mathbf{E}_\pi[\sum_{t \geq 0} (r_{t+1} - J(\theta)) \mid s_0 = s, a_0 = a]$$

- $$\begin{aligned}\nabla_\theta V^\pi(s) &= \nabla_\theta \sum_a \pi(s, a) Q^\pi(s, a) \\ &= \sum_a [\nabla_\theta \pi(s, a) Q^\pi(s, a) + \pi(s, a) \nabla_\theta Q^\pi(s, a)] \\ &= \sum_a \nabla_\theta \pi(s, a) Q^\pi(s, a) + \sum_a \pi(s, a) \nabla_\theta [R_s^a - J(\theta) + \sum_{s'} P_{s,s'}^a V^\pi(s')] \\ &= \sum_a \nabla_\theta \pi(s, a) Q^\pi(s, a) + \sum_{a,s'} \pi(s, a) P_{s,s'}^a \nabla_\theta V^\pi(s') - \nabla_\theta J(\theta)\end{aligned}$$

Proof for average rewards (cont'd)

- Take expectation with respect to the stationary distribution in both sides of the last equation in the previous slide:

$$\begin{aligned}\sum_s \mu^\pi(s) \nabla_\theta V^\pi(s) &= \\&= \sum_{s,a} \mu^\pi(s) \nabla_\theta \pi(s,a) Q^\pi(s,a) + \sum_{s'} \underbrace{\sum_{a,s} \mu^\pi(s) \pi(s,a) P_{s,s'}^a}_{= \mu^\pi(s')} \nabla_\theta V^\pi(s') - \nabla_\theta J(\theta)\end{aligned}$$

- $$\begin{aligned}\Rightarrow \nabla_\theta J(\theta) &= \sum_{s,a} \mu^\pi(s) \nabla_\theta \pi(s,a) Q^\pi(s,a) \\&= \sum_{s,a} \mu^\pi(s) \pi(s,a) \frac{\nabla_\theta \pi(s,a)}{\pi(s,a)} Q^\pi(s,a) \\&= \sum_{s,a} \mu^\pi(s,a) \nabla_\theta \log(\pi(s,a)) Q^\pi(s,a)\end{aligned}$$

Proof for discounted expected rewards

- Basic identities:

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

$$(B) \quad Q^\pi(s, a) = \sum_{s'} P_{s,s'}^a (R_{s,s'}^a + \gamma V^\pi(s'))$$

$$(C) \quad \nabla_\theta V^\pi(s) = \sum_a \nabla_\theta \pi(s, a) Q^\pi(s, a) + \sum_a \pi(s, a) \nabla_\theta Q^\pi(s, a)$$

$$(D) \quad \nabla_\theta Q^\pi(s, a) = \gamma \sum_{s'} P_{s,s'}^a \nabla_\theta V^\pi(s')$$

Proof for discounted expected rewards (cont'd)

$$\bullet \stackrel{(C)}{\nabla_{\theta} V^{\pi}(s)} = \sum_a \nabla_{\theta} \pi(s_0, a) Q^{\pi}(s, a) + \sum_a \pi(s_0, a) \nabla_{\theta} Q^{\pi}(s, a)$$

$$\stackrel{(D)}{=} \sum_a \pi(s, a) \nabla_{\theta} \log(\pi(s, a)) Q^{\pi}(s, a) + \underbrace{\sum_{a,s'} \gamma \pi(s, a) P_{s,s'}^a \nabla_{\theta} V^{\pi}(s')}_S$$

$$\bullet S \stackrel{(C)}{=} \sum_{a,s',a'} \gamma \pi(s, a) P_{s,s'}^a \pi(s', a') \nabla_{\theta} \log(\pi(s', a')) Q^{\pi}(s', a')$$

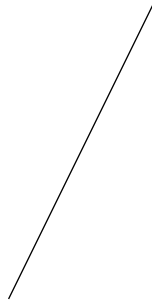
$$+ \sum_{a,s',a'} \gamma \pi(s, a) P_{s,s'}^a \pi(s', a') \nabla_{\theta} Q^{\pi}(s', a')$$

...

Proof for discounted expected rewards (cont'd)

⇒ the policy gradient equation:

$$\nabla_{\theta} V^{\pi}(s) = \sum_{s', a'} \mu^{\pi}(s', a' | s; \gamma) \nabla_{\theta} \log(\pi(s', a')) Q^{\pi}(s', a')$$



$$\mu^{\pi}(s', a' | s; \gamma) := \sum_{t \geq 0} \gamma^t \pi(s', a') \Pr_{\pi}[s_t = s' | s_0 = s]$$

Example 1: softmax policy gradient

- State-action pairs weighted by linear combination of features:

$$\pi(s, a; \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\phi}(s,a)^\top \boldsymbol{\theta}}}{\sum_{a'} e^{\boldsymbol{\phi}(s,a')^\top \boldsymbol{\theta}}}$$

- The score function:

$$\frac{\partial}{\partial \theta_i} \log(\pi(s, a; \boldsymbol{\theta})) = \phi_i(s, a) - \frac{\sum_{a'} \phi_i(s, a') e^{\boldsymbol{\phi}(s,a')^\top \boldsymbol{\theta}}}{\sum_{a'} e^{\boldsymbol{\phi}(s,a')^\top \boldsymbol{\theta}}}$$

or, equivalently,

$$\nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta})) = \boldsymbol{\phi}(s, a) - \mathbf{E}_{a' \sim \pi(s, \cdot; \boldsymbol{\theta})}[\boldsymbol{\phi}(s, a')]$$

Example 2: continuous action space

- Actions state space: set of real numbers $A = \mathbf{R}$
- Policy approximator:

$$\pi(s, a; \theta) = \frac{1}{\sqrt{2\pi}\sigma(s; \theta)} \exp\left(-\frac{(a - \mu(s; \theta))^2}{2\sigma(s; \theta)^2}\right)$$

where $\mu(s; \theta)$ and $\sigma(s; \theta)$ are mean and deviation function approximators

- Linear function approximators with feature vectors $\phi_\mu(s)$ and $\phi_\sigma(s)$:

- $\mu(s; \theta) = \phi_\mu(s)^\top \theta_\mu$ and $\sigma(s; \theta) = \phi_\sigma(s)^\top \theta_\sigma$

- $\nabla_{\theta_\mu} \log(\pi(s, a; \theta)) = \frac{a - \mu(s; \theta)}{\sigma(s; \theta)^2} \phi_\mu(s)$

- $\nabla_{\theta_\sigma} \log(\pi(s, a; \theta)) = \frac{(a - \mu(s; \theta))^2 - \sigma(s; \theta)^2}{\sigma(s; \theta)^2} \phi_\sigma(s)$

(Exercise: check this)

Example 3: Bernoulli, logistic example

- Action state space: $A = \{0,1\}$
- Policy approximator:

$$\pi(a = 1, s; \theta) = 1 - \pi(a = 0, s; \theta) := p(s; \theta)$$

where $p(s; \theta)$ is a function approximator

- Linear function approximator with feature vectors $\phi(s)$
 - Each state has action preference scores $h_0(s; \theta)$ and $h_1(s; \theta)$ such that

$$h_1(s; \theta) - h_0(s; \theta) = \phi(s)^\top \theta$$

- **Fact 1:** for exponential soft-max policy $p(s; \theta) = \sigma(\phi(s)^\top \theta)$ logistic function

- **Fact 2:** $\nabla_\theta \log(\pi(s, a; \theta)) = (a - \sigma(\phi(s)^\top \theta))\phi(s)$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

REINFORCE: MC policy gradient algorithm

- Use return R_t as an estimate of $Q^{\pi(\cdot; \theta)}(s_t, a_t)$
- **Initialization:** θ arbitrary
- **For each** episode $(s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T)$ generated using policy $\pi(\cdot; \theta)$:

For $t = 1, 2, \dots, T - 1$ **do**:

$$\theta \leftarrow \theta + \eta \nabla_{\theta} \log(\pi(s_t, a_t; \theta)) R_t$$

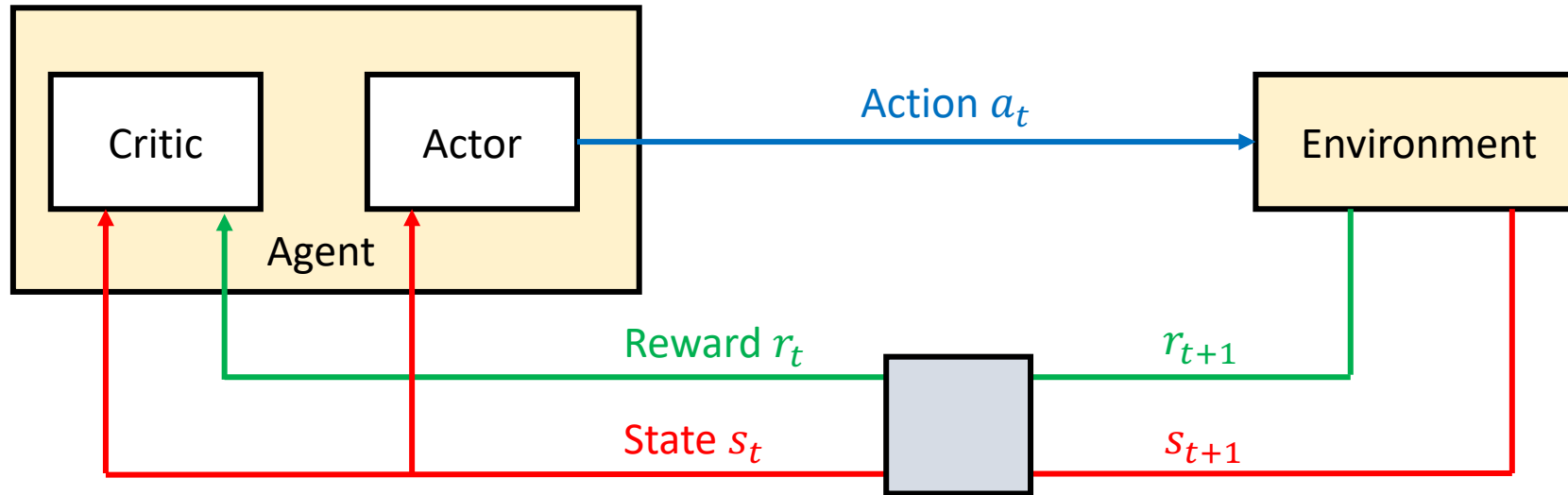
end for

return θ

Actor-critic algorithms

- MC policy gradient algorithm may have a high variance
- Solution sought by using actor-critic algorithms
- Actor-critic algorithms combine policy gradient with value function estimation

Actor-critic algorithms



- Critic uses an approximator to learn a value function
- Actor updates policy approximator in a direction of performance improvement

Actor-critic control

- **Critic:** estimates $Q^{\pi(\cdot; \theta)}(s, a)$ by an approximator $\hat{Q}(s, a; \mathbf{w})$
 - The critic performs policy evaluation
 - Standard methods can be applied: MC, temporal difference learning, TD(λ), gradient-based least-square estimation methods, ...
- **Actor:** updates policy parameter θ in direction suggested by the critic
 - The actor performs control using approximate policy gradient:

$$\nabla_{\theta} J(\theta) \approx \mathbf{E}[\nabla_{\theta} \log(\pi(s, a; \theta)) \hat{Q}(s, a; \mathbf{w})]$$

- Parameter update:

$$\theta \leftarrow \theta + \eta \nabla_{\theta} \log(\pi(s, a; \theta)) \hat{Q}(s, a; \mathbf{w})$$

Example actor-critic algorithm

- **Initialization:** s, θ, w
- **For each** episode:
 Sample action a from $\pi(s, \cdot; \theta)$
 Repeat until s is terminal:

 Receive reward r and next state s'
 Sample action a' from $\pi(s', \cdot; \theta)$

$$\delta \leftarrow r + \gamma \hat{Q}(s', a'; w) - \hat{Q}(s, a; w)$$

$$\theta \leftarrow \theta + \eta \nabla_{\theta} \log(\pi(s, a; \theta)) \hat{Q}(s, a; w)$$

$$w \leftarrow w + \alpha \delta \phi(s, a)$$

$$a \leftarrow a', s \leftarrow s'$$

Linear value function approximator:

$$\hat{Q}(s, a; w) = \phi(s, a)^{\top} w$$

Critic: updates w by linear TD(0)

Actor: updates θ by policy gradient

Bias issue of approximate policy gradient

- Using approximate policy gradient introduces a bias
- A biased policy gradient may fail to converge to a globally optimal solution
- Unbiased policy gradient can be ensured under certain conditions (next slide)

Compatible function approximation theorem

- **Thm.** Assume the following two conditions:
 - (C1) compatibility of value function approximator and the policy:

$$\nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) = \nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta}))$$

- (C2) value function approximator minimizes the mean-squared error:

$$\text{MSE}(\mathbf{w}) = \mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} \left[\left(Q^{\pi(\cdot; \boldsymbol{\theta})}(s, a) - \hat{Q}(s, a; \mathbf{w}) \right)^2 \right]$$

- Then, the policy gradient is exact:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta})) \hat{Q}(s, a; \mathbf{w})]$$

Proof sketch

- First, note that by (C2), $\nabla_{\mathbf{w}} \text{MSE}(\mathbf{w}) = 0$, i.e.

$$\mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} \left[\left(Q^{\pi(\cdot; \boldsymbol{\theta})}(s, a) - \hat{Q}(s, a; \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) \right] = 0$$

- By (C1),

$$\mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} \left[\left(Q^{\pi(\cdot; \boldsymbol{\theta})}(s, a) - \hat{Q}(s, a; \mathbf{w}) \right) \nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta})) \right] = 0$$

which is equivalent to

$$\mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} \left[\hat{Q}(s, a; \mathbf{w}) \nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta})) \right] = \mathbf{E}_{\pi(\cdot; \boldsymbol{\theta})} \left[\underbrace{Q^{\pi(\cdot; \boldsymbol{\theta})}(s, a) \nabla_{\boldsymbol{\theta}} \log(\pi(s, a; \boldsymbol{\theta}))}_{\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})} \right]$$

Compatible function approximation in action

- Consider the soft-max policy, for given state-action feature vectors $\phi(s, a)$:

$$\pi(s, a; \theta) = \frac{e^{\phi(s, a)^\top \theta}}{\sum_{a'} e^{\phi(s, a')^\top \theta}}$$

- Compatibility condition requires that

$$\nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) = \nabla_{\theta} \log(\pi(s, a; \theta)) = \phi(s, a) - \sum_{a'} \phi(s, a') \pi(s, a'; \theta)$$

which leads to a linear approximator for the value function:

$$\hat{Q}(s, a; \mathbf{w}) = \underbrace{(\phi(s, a) - \sum_{a'} \phi(s, a') \pi(s, a'; \theta))^\top}_{\text{centered state-action feature vectors}} \mathbf{w}$$

Convergence theorem

- **Thm.** Assume:
 - $\pi(\cdot; \theta)$ and $\hat{Q}(\cdot; \mathbf{w})$ are differentiable functions
 - Compatibility condition holds
 - The Hessian matrix $\nabla_{\theta}^2 \pi(s, a; \theta)$ satisfies $\|\nabla_{\theta}^2 \pi(s, a; \theta)\|_{\infty} \leq B < \infty$
 - Step sizes are such that $\lim_{t \rightarrow \infty} \eta_t = 0$ and $\sum_{t > 0} \eta_t = \infty$
- Then, for any MDP with bounded rewards and the parameter updates:

$$\theta_{t+1} = \theta_t + \eta_t \sum_{s,a} \mu^{\pi(\cdot; \theta_t)}(s) \pi(s, a; \theta_t) \nabla_{\theta} \log(\pi(s, a; \theta_t)) \hat{Q}(s, a; \mathbf{w}_t)$$

where \mathbf{w}_t is a solution of the equation:

$$\sum_{s,a} \mu^{\pi(\cdot; \theta_t)}(s) \pi(s, a; \theta_t) [Q^{\pi(\cdot; \theta_t)}(s, a) - \hat{Q}(s, a; \mathbf{w})] \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w}) = 0$$

are convergent in the sense that $\lim_{t \rightarrow \infty} \|\nabla_{\theta} J(\theta_t)\| = 0$

Separation of timescales

- In the last theorem, \mathbf{w}_t is defined as a solution of a fixed point equation which has the policy's parameter vector $\boldsymbol{\theta}_t$ as a parameter
- An incremental implementation would update \mathbf{w}_t incrementally in a similar manner as for the parameter vector $\boldsymbol{\theta}_t$ but with a larger step size parameter
- This can be seen as a **separation of timescales**:
 - Critic updates the value function approximator at a faster timescale trying to evaluate the current policy chosen by the actor
 - Actor varies the policy's parameter more slowly to allow the critic to evaluate the current policy
- Separation of timescales assumption is common in convergence proofs of actor-critic algorithms

Variance reduction using a baseline

- Subtracting a baseline function $B(s)$ from the policy gradient

$$\nabla_{\theta} \log(\pi(s, a; \theta)) (Q^{\pi(\cdot; \theta)}(s, a) - B(s))$$

does not change the expected policy gradient:

$$\begin{aligned} \mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) (Q^{\pi(\cdot; \theta)}(s, a) - B(s))] \\ = \mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) Q^{\pi(\cdot; \theta)}(s, a)] \end{aligned}$$

which is because

$$\mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) B(s)] = 0 \quad (\text{check this})$$

Policy gradient using the advantage function

- Choose the baseline $B(s) = V^{\pi(\cdot; \theta)}(s)$
- The advantage function: $A^{\pi(\cdot; \theta)}(s, a) := Q^{\pi(\cdot; \theta)}(s, a) - V^{\pi(\cdot; \theta)}(s)$
- Policy gradient using the advantage function:

$$\nabla_{\theta} J(\theta) = \mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) A^{\pi(\cdot; \theta)}(s, a)]$$

- The advantage function can reduce variance of policy gradient

An approach for estimating the advantage function

- The critic may compute estimators of both value functions:

$$\hat{Q}(s, a; \mathbf{w}) \text{ for } Q^{\pi(\cdot; \boldsymbol{\theta})}(s, a)$$

and

$$\hat{V}(s, a; \mathbf{v}) \text{ for } V^{\pi(\cdot; \boldsymbol{\theta})}(s)$$

which can be done by standard methods such as TD learning

- The estimator of the advantage function:

$$\hat{A}(s, a) = \hat{Q}(s, a; \mathbf{w}) - \hat{V}(s, a; \mathbf{v})$$

Another approach for estimating the advantage function

- The TD error $\delta^{\pi(\cdot;\theta)}(s, s', r) = r + \gamma V^{\pi(\cdot;\theta)}(s') - V^{\pi(\cdot;\theta)}(s)$ is an **unbiased estimate** of the advantage function:

$$\begin{aligned}\mathbf{E}_{\pi(\cdot;\theta)}[\delta^{\pi(\cdot;\theta)}(s_t, s_{t+1}, r_{t+1}) | s_t = s, a_t = a] \\&= \mathbf{E}_{\pi(\cdot;\theta)}[r_{t+1} + \gamma V^{\pi(\cdot;\theta)}(s_{t+1}) | s_t = s, a_t = a] - V^{\pi(\cdot;\theta)}(s) \\&= Q^{\pi(\cdot;\theta)}(s, a) - V^{\pi(\cdot;\theta)}(s)\end{aligned}$$

⇒ we may use the TD error in the policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbf{E}_{\pi(\cdot;\theta)}[\nabla_{\theta} \log(\pi(s, a; \theta)) \delta^{\pi(\cdot;\theta)}]$$

- In practice, we will use an approximate TD error:

$$\delta^{\pi(\cdot;\theta)}(s, s', r) = r + \gamma \hat{V}(s'; \mathbf{v}) - \hat{V}(s; \mathbf{v})$$

Critic policy evaluation methods

- The critic can use different targets to evaluate:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left(v_t - \hat{V}(s; \mathbf{w}) \right) \phi(s)$$

- The target v_t is defined differently for different methods:

- MC:

$$v_t = R_t$$

- TD(0):

$$v_t = r_{t+1} + \gamma \hat{V}(s_{t+1}; \mathbf{w})$$

- Forward-view TD(λ):

$$v_t = R_t^\lambda \quad (\lambda\text{-return})$$

- Backward-view TD(λ):

$$\delta_t = r_{t+1} + \gamma \hat{V}(s_{t+1}; \mathbf{w}) - \hat{V}(s_t; \mathbf{w})$$

$$e_t = \gamma \lambda e_{t-1} + \phi(s_t)$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta_t e_t$$

Actor policy gradient methods

- The policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) A^{\pi(\cdot; \theta)}(s, a)]$$

- Gradient ascent method:

$$\theta \leftarrow \theta + \eta \nabla_{\theta} \log(\pi(s, a; \theta)) \hat{A}(s, a; \nu)$$

- Examples:

- MC: $\hat{A}(s, a; \mathbf{w}) = R_t - \hat{V}(s; \mathbf{w})$
- One-step TD error: $\hat{A}(s, a; \mathbf{w}) = r + \gamma \hat{V}(s'; \mathbf{w}) - \hat{V}(s; \mathbf{w})$

Summary

- Equivalent form of policy gradient: $\nabla_{\theta} J(\theta) = \mathbf{E}_{\pi(\cdot; \theta)} [\nabla_{\theta} \log(\pi(s, a; \theta)) v]$

Method	v
MC (REINFORCE)	R_s^a
Q actor-critic	$\hat{Q}(s, a; w)$
Advantage actor-critic	$\hat{A}(s, a; w)$
TD actor-critic	$\delta^{\pi(\cdot; \theta)}$
TD(λ) actor-critic	$\delta^{\pi(\cdot; \theta)} e^{\pi(\cdot; \theta)}$

- Use stochastic gradient descent algorithm
- Critic variants that use policy evaluation to estimate $Q^{\pi(\cdot; \theta)}$, $V^{\pi(\cdot; \theta)}$, or $A^{\pi(\cdot; \theta)}$

References

- Sutton and Barto, Reinforcement Learning, The MIT Press, 2nd edition, 2018, Chapter 13: Policy gradient methods
- Sutton et al, Policy Gradient Methods for Reinforcement Learning with Function Approximation, NIPS 2000
- Silver, Lecture 7: [Policy Gradient](#)

Seminar exercises

- Solving Atari game Breakout using policy gradient method

Extras

Reshaping rewards

- Consider an MDP(P, R) with transition probabilities and expected rewards

$$P = (P_{s,s'}^a) \text{ and } R = (R_{s,s'}^a)$$

- Consider a transformation of expected rewards as follows

$$\tilde{R}_{s,s'}^a = R_{s,s'}^a + \gamma B(s') - B(s)$$

for an arbitrarily given function $B: S \rightarrow \mathbf{R}$

- For example, the transformation of expected rewards holds if instantaneous rewards are transformed as $\tilde{r}_{t+1} = r_{t+1} + \gamma B(s_{t+1}) - B(s_t)$

Reshaping rewards theorem

- **Thm**: MDP(P, R) and MDP(P, \tilde{R}) have the same optimal policies.
- Exercise: show this

Proof sketch

- Let Q^* be the optimal action value function and π^* be an optimal policy for $\text{MDP}(P, R)$, and let \tilde{Q}^* be the optimal value function for $\text{MDP}(P, \tilde{R})$
- By the Bellman's optimality equation:

$$Q^*(s, a) = \sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma \max_{a'} Q^*(s', a')] \text{ for all state action pairs } (s, a)$$

- Since $\sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma \max_{a'} Q^*(s', a')]$
$$= \sum_{s'} P_{s,s'}^a [\tilde{R}_{s,s'}^a + B(s) + \gamma \max_{a'} [Q^*(s', a') - B(s)]]$$

we have $\underbrace{[Q^*(s, a) - B(s)]}_{\tilde{Q}^*(s, a)} = \sum_{s'} P_{s,s'}^a [\tilde{R}_{s,s'}^a + \gamma \max_{a'} \underbrace{[Q^*(s', a') - B(s)]}_{\tilde{Q}^*(s', a')}]$