

# Planning and Optimal Control

## 8. Policy Gradient Methods

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)  
Institute for Computer Science  
University of Hildesheim, Germany

# Syllabus

## A. Models for Sequential Data

- Tue. 22.10. (1) 1. Markov Models
- Tue. 29.10. (2) 2. Hidden Markov Models
- Tue. 5.11. (3) 3. State Space Models
- Tue. 12.11. (4) 3b. (ctd.)

## B. Models for Sequential Decisions

- Tue. 19.11. (5) 1. Markov Decision Processes
- Tue. 26.11. (6) 1b. (ctd.)
- Tue. 3.12. (7) 1c. (ctd.)
- Tue. 10.12. (8) 2. Monte Carlo and Temporal Difference Methods
- Tue. 17.12. (9) 3. Q Learning
- Tue. 24.12. — — *Christmas Break* —
- Tue. 7.1. (10) 4. Policy Gradient Methods
- Tue. 14.1. (11) tba
- Tue. 21.1. (12) tba
- Tue. 28.1. (13) 8. Reinforcement Learning for Games
- Tue. 4.2. (14) Q&A

# Outline

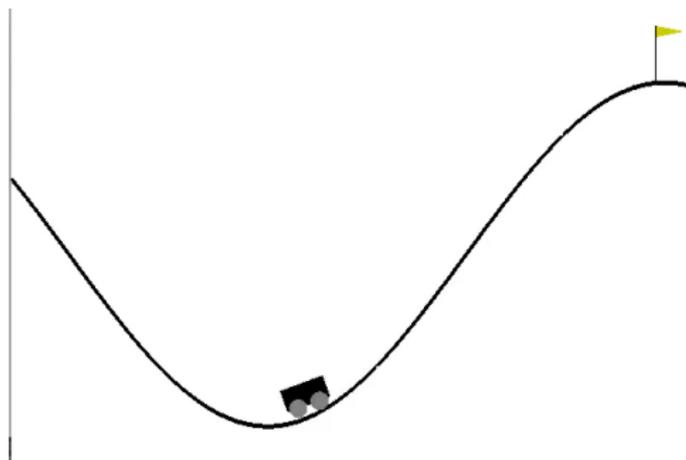
1. The Policy Gradient Theorem
2. Monte Carlo Policy Gradient (REINFORCE)
3. TD Policy Gradients: Actor-Critic Methods
4. Deterministic Policy Gradients

# Outline

1. The Policy Gradient Theorem
2. Monte Carlo Policy Gradient (REINFORCE)
3. TD Policy Gradients: Actor-Critic Methods
4. Deterministic Policy Gradients

- ▶ Q Learning: learn the action value function.
- ▶ Policy Gradient Methods

# Example: Continuous Mountain Car



$S := [-1.2, 0.5] \times [-0.07, 0.07]$ ,  $s := (x, v)$  (position and velocity)

$A := [-1, +1]$  (acceleration)

$x_{t+1} := \text{clip}(x_t + v_t)$

$v_{t+1} := \text{clip}(v_t + 0.001a_t - 0.0025 \cos(3x_t))$

$p(x_0) := \text{unif}([-0.6, -0.4])$ ,  $v_0 := 0$

# Idea

- ▶ let  $\pi$  be a parametrized policy with parameters  $\theta$ :

$$\pi(a \mid s; \theta)$$

- ▶ i.e., a neural network with input  $s$  and output  $a$
- ▶ view its value function  $V^\pi$  of the unique starting state  $s_0$  as a function of  $\theta$ :

$$V(\theta) := V^\pi(s_0; \theta)$$

- ▶ finding the optimal policy  $\rightsquigarrow$  maximize  $V$  w.r.t.  $\theta$ 
  - ▶ e.g., by gradient ascent:

$$\theta^{(t+1)} := \theta^{(t)} + \alpha_t \nabla_\theta V(\theta^{(t)})$$

# State Distribution of a Policy

state visiting frequencies of a policy  $\pi$ :

$$\eta^\pi(s) := \mathbb{E}(|\{S_t = s \mid t = 1 : T\}|)$$

consistency:

$$\eta^\pi(s) = p(S_0 = s) + \sum_{s' \in \mathcal{S}} \eta^\pi(s') \sum_{a \in A} \pi(a \mid s') p(s \mid s', a)$$

$$\rightsquigarrow \eta^\pi = (I - ((P^\pi)^T)^{-1}) p_0$$

with  $P^\pi := (\sum_{a \in A} \pi(a \mid s') p(s \mid s', a))_{(s', s) \in \mathcal{S}^2}$  state transition under  $\pi$

$p_0 := (p(S_0 = s))_{s \in \mathcal{S}}$  initial states

state distribution of policy  $\pi$  (**on-policy distribution**):

$$\mu^\pi(s) := \frac{\eta^\pi(s)}{\sum_{s' \in \mathcal{S}} \eta^\pi(s')}$$

Note: With  $I$  the identity matrix.

## Discounted State Distribution of a Policy

**discounted** state visiting frequencies of a policy  $\pi$ :

$$\eta^\pi(s) := \mathbb{E}\left(\sum_t \gamma^t \mathbb{I}(S_t = s)\right)$$

- ▶ visiting a state at time  $t$  contributes weight  $\gamma^t$ .
- ▶ = visiting frequency for  $\gamma = 1$ .

consistency:

$$\eta^\pi(s) = p(S_0 = s) + \sum_{s' \in \mathcal{S}} \eta^\pi(s') \gamma \sum_{a \in A} \pi(a | s') p(s | s', a)$$

$$\rightsquigarrow \eta^\pi = (I - \gamma((P^\pi)^T)^{-1})p_0$$

**discounted** state distribution of policy  $\pi$ :

$$\mu^\pi(s) := \frac{\eta^\pi(s)}{\sum_{s' \in \mathcal{S}} \eta^\pi(s')}$$

# Policy Gradient Theorem

$$\nabla_{\theta} V^{\pi}(s_0) \propto \sum_s \mu^{\pi}(s) \sum_a \nabla \pi(a | s; \theta) Q^{\pi}(s, a)$$

# Policy Gradient Theorem / Proof

$$\begin{aligned}
 \nabla V^\pi(s) &= \nabla \sum_a \pi(a | s) Q^\pi(s, a) \\
 &= \sum_a \nabla \pi(a | s) Q^\pi(s, a) + \pi(a | s) \nabla Q^\pi(s, a) \\
 &= \sum_a \nabla \pi(a | s) Q^\pi(s, a) + \pi(a | s) \nabla \sum_{s', r} p(s', r | s, a) (r + \gamma V^\pi(s')) \\
 &= \sum_a \nabla \pi(a | s) Q^\pi(s, a) + \pi(a | s) \sum_{s', r} p(s' | s, a) \gamma \nabla V^\pi(s') \\
 &= \text{rec.} \sum_a \nabla \pi(a | s) Q^\pi(s, a) + \pi(a | s) \sum_{s'} p(s' | s, a) \gamma \\
 &\quad \left( \sum_{a'} \nabla \pi(a' | s') Q^\pi(s', a') + \pi(a' | s') \sum_{s''} p(s'' | s', a') \gamma \nabla V^\pi(s'') \right) \\
 &= \text{rec.} \sum_{s'} \sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \gamma^k \sum_a \nabla \pi(a | s') Q^\pi(s', a)
 \end{aligned}$$

Note:  $\gamma$  is missing in [Sutton and Barto, 2018, p.325].

# Policy Gradient Theorem / Proof

$$\nabla V^\pi(s) = \sum_{s'} \sum_{k=0}^{\infty} \Pr(s \rightarrow s', k, \pi) \gamma^k \sum_a \nabla \pi(a | s') Q^\pi(s', a)$$

$$\begin{aligned} \nabla V^\pi(s_0) &= \sum_s \sum_{k=0}^{\infty} \Pr(s_0 \rightarrow s, k, \pi) \gamma^k \sum_a \nabla \pi(a | s) Q^\pi(s, a) \\ &= \sum_s \eta^\pi(s) \sum_a \nabla \pi(a | s) Q^\pi(s, a) \\ &\propto \sum_s \mu^\pi(s) \sum_a \nabla \pi(a | s) Q^\pi(s, a) \end{aligned}$$

# Outline

1. The Policy Gradient Theorem
2. Monte Carlo Policy Gradient (REINFORCE)
3. TD Policy Gradients: Actor-Critic Methods
4. Deterministic Policy Gradients

# Monte Carlo Policy Gradient (REINFORCE)

$$\begin{aligned}
 \nabla V^\pi(s_0) &\propto \sum_s \mu^\pi(s) \sum_a \nabla \pi(a | s) Q^\pi(s, a) \\
 &= \mathbb{E} \left( \sum_a \nabla \pi(a | S_t) Q^\pi(S_t, a) \right) \\
 &= \mathbb{E} \left( \sum_a \frac{\nabla \pi(a | S_t)}{\pi(a | S_t)} \pi(a | S_t) Q^\pi(S_t, a) \right) \\
 &= \mathbb{E} \left( \frac{\nabla \pi(A_t | S_t)}{\pi(A_t | S_t)} V_t \right) \\
 &= \mathbb{E} (V_t \nabla \log \pi(A_t | S_t))
 \end{aligned}$$

↪ update rule:

$$\theta := \theta + \alpha v_t \nabla \log \pi(a_t | s_t)$$

Note: Observed value  $V_t$  also called **return** in the literature.

# Monte Carlo Policy Gradient (REINFORCE) / Alg.

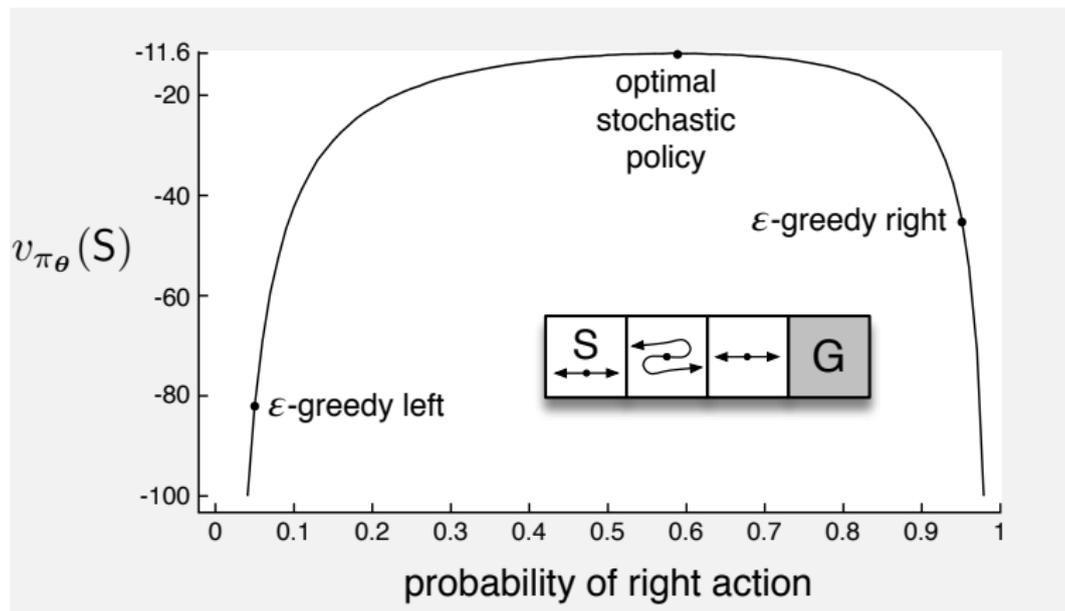
```

1 learn-opt-policy-discounted-mc-policygrad( $S, A, \gamma, s_{\text{term}}, N, \alpha$ ):
2   initialize parameters  $\theta$  of policy  $\pi$ 
3   for  $n := 1, \dots, N$ :
4     ( $s, a, r, T$ ) := generate-episode( $S, A, s_{\text{term}}, \pi$ )
5     for  $t := 0, \dots, T - 1$ :
6        $v_t := \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$ 
7        $\theta := \theta + \alpha v_t \nabla \log \pi(a_t | s_t)$ 
8   return  $\pi$ 
  
```

where

- ▶  $s_{\text{term}}$  terminal state with zero reward.
- ▶  $\alpha$  **learning rate**
- ▶ returning  $\pi$  actually means to return its parameters  $\theta$

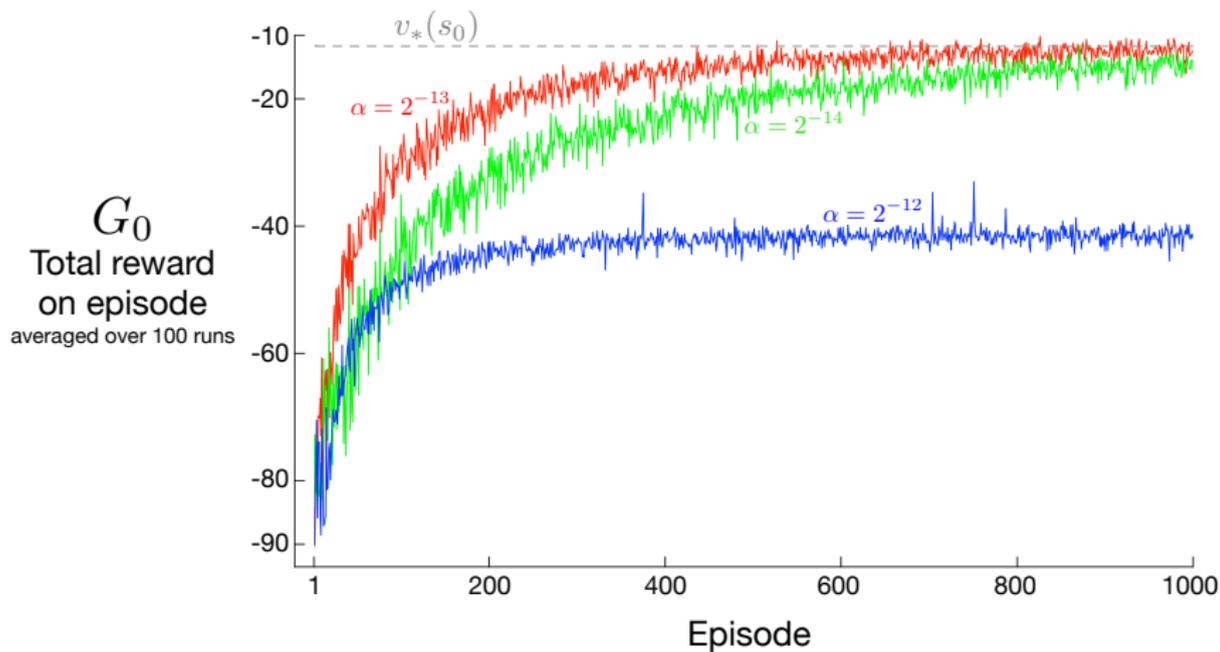
# Example



[source: [Sutton and Barto, 2018, p.323]]

- ▶ assume all states are indistinguishable; reward  $-1$  per step
- ▶  $\epsilon := 0.1$

# Example



[source: [Sutton and Barto, 2018, p.328]]

# Monte Carlo Policy Gradient (REINFORCE w. Baseline)

- **baseline**  $b : S \rightarrow \mathbb{R}$ , not depending on actions  $a$ .

$$\begin{aligned} \nabla V^\pi(s_0) &\propto \sum_s \mu^\pi(s) \sum_a \nabla \pi(a | s) Q^\pi(s, a) \\ &= \sum_s \mu^\pi(s) \sum_a \nabla \pi(a | s) (Q^\pi(s, a) - b(s)) \end{aligned}$$

as

$$\sum_a \nabla \pi(a | s) b(s) = b(s) \nabla \sum_a \pi(a | s) = b(s) \nabla 1 = 0$$

# Monte Carlo Policy Gradient (REINFORCE w. Baseline)

$$\begin{aligned}
 \nabla V^\pi(s_0) &\propto \sum_s \mu^\pi(s) \sum_a \nabla \pi(a | s) (Q^\pi(s, a) - b(s)) \\
 &= \mathbb{E} \left( \sum_a \nabla \pi(a | S_t) (Q^\pi(S_t, a) - b(S_t)) \right) \\
 &= \mathbb{E} \left( \sum_a \frac{\nabla \pi(a | S_t)}{\pi(a | S_t)} \pi(a | S_t) (Q^\pi(S_t, a) - b(S_t)) \right) \\
 &= \mathbb{E} \left( \frac{\nabla \pi(A_t | S_t)}{\pi(A_t | S_t)} (V_t - b(S_t)) \right) \\
 &= \mathbb{E} \left( (V_t - b(S_t)) \nabla \log \pi(A_t | S_t) \right)
 \end{aligned}$$

↪ update rule:

$$\theta := \theta + \alpha (v_t - b(s_t)) \nabla \log \pi(a_t | s_t)$$

# Monte Carlo Policy Gradient (REINFORCE w. Baseline)

- ▶ often a current estimate of the value function  $\hat{V}$  is used as baseline:

$$b(s) := \hat{V}(s; \eta)$$

- ▶ updates for parameters  $\eta$  of  $\hat{V}$ :

$$\eta := \eta + \alpha_2 (v_t - \hat{V}(s_t)) \nabla \hat{V}(s_t)$$

# Monte Carlo Policy Gradient (REINFORCE w. Baseline)

## Alg.

1 **learn-opt-policy-discounted-mc-policygrad-base**( $S, A, \gamma, s_{\text{term}}, N, \alpha$ ):

2 **initialize parameters**  $\eta$  of value function  $\hat{V}$

3 initialize parameters  $\theta$  of policy  $\pi$

4 for  $n := 1, \dots, N$ :

5  $(s, a, r, T) := \text{generate-episode}(S, A, s_{\text{term}}, \pi)$

6 for  $t := 0, \dots, T - 1$ :

7  $v_t := \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$

8  $\hat{v}_t := \hat{V}(s_t)$

9  $\eta := \eta + \alpha_2 (v_t - \hat{v}_t) \nabla V(s_t)$

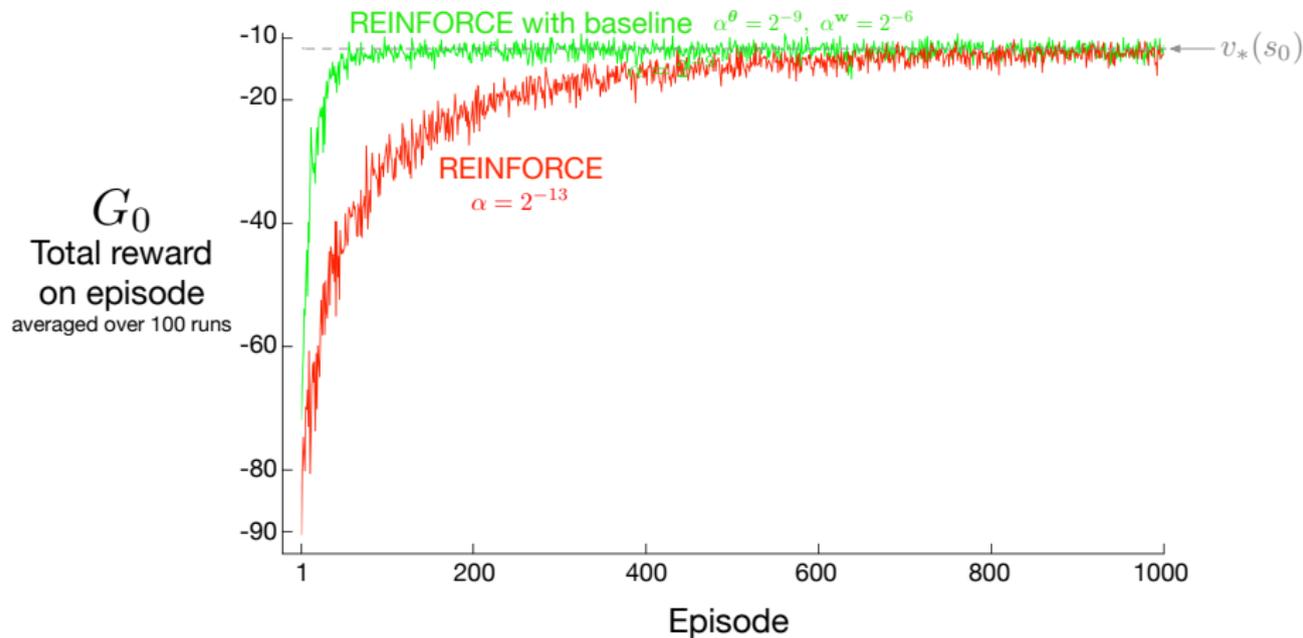
10  $\theta := \theta + \alpha_1 (v_t - \hat{v}_t) \nabla \log \pi(a_t | s_t)$

11 return  $\pi$

where

- ▶  $s_{\text{term}}$  terminal state with zero reward.
- ▶  $\alpha_1, \alpha_2$  **learning rates** for  $\pi$  and  $\hat{V}$ , resp.
- ▶ returning  $\pi$  actually means to return its parameters  $\theta$

# Example



[source: [Sutton and Barto, 2018, p.330]]

# Outline

1. The Policy Gradient Theorem
2. Monte Carlo Policy Gradient (REINFORCE)
3. TD Policy Gradients: Actor-Critic Methods
4. Deterministic Policy Gradients

# TD Policy Gradients: Actor-Critic Methods

To derive MC policy gradient / REINFORCE:

$$\mathbb{E}\left(\sum_a \pi(a | S_t) Q^\pi(S_t, a)\right) = \mathbb{E}(V_t)$$

To derive TD policy gradient / Actor Critic:

$$\mathbb{E}\left(\sum_a \pi(a | S_t) Q^\pi(S_t, a)\right) = \mathbb{E}(R_t + \gamma V_{t+1})$$

then plug in  $\hat{V}(S_{t+1})$  for  $V_{t+1}$ :

$$\nabla V^\pi(s_0) \underset{\text{approx.}}{\propto} \mathbb{E}((R_t + \gamma \hat{V}(S_{t+1})) \nabla \log \pi(A_t | S_t))$$

and subtract baseline  $\hat{V}(S_t)$ :

$$\nabla V^\pi(s_0) \underset{\text{approx.}}{\propto} \mathbb{E}((R_t + \gamma \hat{V}(S_{t+1}) - \hat{V}(S_t)) \nabla \log \pi(A_t | S_t))$$

$\rightsquigarrow$  update rule:

$$\theta := \theta + \alpha (r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)) \nabla \log \pi(a_t | s_t)$$

# TD Policy Gradients (Actor Critic)

```

1 learn-opt-policy-discounted-td-policygrad( $S, A, \gamma, s_{\text{term}}, N, \alpha$ ):
2   initialize parameters  $\eta$  of value function  $\hat{V}$ 
3   initialize parameters  $\theta$  of policy  $\pi$ 
4   for  $n := 1, \dots, N$ :
5      $s := \text{new\_process}()$ 
6      $\hat{v} := \hat{V}(s)$ 
7     while  $s \neq s_{\text{term}}$ :
8        $a := \pi(s)$ 
9        $(r, s') := \text{execute\_action}(s, a)$ 
10       $\hat{v}' := \hat{V}(s')$ 
11       $\delta := r + \gamma \hat{v}' - \hat{v}$ 
12       $\eta := \eta + \alpha_2 \delta \nabla V(s)$ 
13       $\theta := \theta + \alpha_1 \delta \nabla \log \pi(a | s)$ 
14       $s := s', \hat{v} := \hat{v}'$ 
15   return  $\pi$ 
  
```

where

- ▶ `new_process()` sets up a new process.
- ▶ `execute_action(s, a)` executes action  $a$  in process in state  $s$ .

limitations:

- ▶ online  $\rightsquigarrow$  replay memory
- ▶ gradient descent  $\rightsquigarrow$  general model update algorithm

# TD Policy Gradients (Actor Critic) / Replay Memory

```

1 learn-opt-policy-discounted-td-policygrad( $S, A, \gamma, s_{\text{term}}, N, \alpha$ ):
2   initialize parameters  $\eta$  of value function  $\hat{V}$ 
3   initialize parameters  $\theta$  of policy  $\pi$ 
4    $\mathcal{D} := \emptyset$ 
5   for  $n := 1, \dots, N$ :
6      $s := \text{new\_process}()$ 
7     while  $s \neq s_{\text{term}}$ :
8        $a := \pi(s)$ 
9        $(r, s') := \text{execute\_action}(s, a)$ 
10       $\mathcal{D} := \mathcal{D} \cup \{(s, a, r, s')\}$ 
11       $\mathcal{D}'_1 := \{(s, a, r + \gamma \hat{V}(s')) \mid (s, a, r, s') \in \mathcal{D}\}$ 
12       $\mathcal{D}'_2 := \{(s, a, \text{caseweight} = r + \gamma \hat{V}(s') - \hat{V}(s)) \mid (s, a, r, s') \in \mathcal{D}\}$ 
13       $\hat{V} := \text{update-model}(\hat{V}, \mathcal{D}'_1)$ 
14       $\pi := \text{update-model}(\pi, \mathcal{D}'_2)$ 
15       $s := s'$ 
16   return  $\pi$ 
  
```

where

- ▶ `new_process()` sets up a new process.
- ▶ `execute_action(s, a)` executes action  $a$  in process in state  $s$ .

# Outline

1. The Policy Gradient Theorem
2. Monte Carlo Policy Gradient (REINFORCE)
3. TD Policy Gradients: Actor-Critic Methods
4. Deterministic Policy Gradients

# Policy Value

- ▶ for non-finite MDPs, use expected value over starting states (policy value) to define optimality of a policy:

$$\begin{aligned}
 V(\pi) &:= \mathbb{E}_{s_0 \sim p}(V^\pi(s_0)) = \int_{s_0} V^\pi(s_0) p(s_0) ds_0 \\
 &= \mathbb{E}_{s \sim \mu^\pi, a \sim \pi}(r(s, a))
 \end{aligned}$$

- ▶ policy gradient theorem is for **stochastic policies**

$$\pi : S \times A \rightarrow [0, 1]$$

- ▶ is there a corresponding result for **deterministic policies**?

$$\pi : S \rightarrow A$$

# Deterministic Policy Gradients

- ▶ stochastic policy gradient:

$$\nabla V(\pi) = \mathbb{E}_{s \sim \eta^\pi, a \sim \pi} (\nabla_\theta \log \pi(s, a; \theta) Q^\pi(s, a))$$

- ▶ deterministic policy gradient:

$$\nabla V(\pi) = \mathbb{E}_{s \sim \eta^\pi} (\nabla_\theta \pi(s; \theta) \nabla_a Q^\pi(s, a) |_{a=\pi(s)})$$

- ▶ deterministic policy gradient is limit of stochastic one:

$$\lim_{\pi' \rightarrow \pi} \nabla_\theta V(\pi') = \nabla_\theta V(\pi), \quad \pi' \text{ stochastic policy, } \pi \text{ deterministic policy}$$

- ▶ see Silver et al. [2014]

# Overview

		<i>from what to learn</i>	
		from episodes	from transitions
<i>what to learn</i>	value function $V^\pi$	Monte Carlo for $V^\pi$	Temporal Differences TD
	action value function $Q^\pi$	Monte Carlo for $Q^\pi$	SARSA
	optimal action value function $Q^{\pi^*}$	.	Q Learning
	optimal policy $\pi^*$	Monte Carlo Policy Gradient / REINFORCE	TD Policy Gradient / Actor Critic

# Summary

- ▶ **Policy gradient methods** parametrize the policy directly:  $\pi(s; \theta)$ 
  - ▶ instead of parametrizing the action-value function  $Q(s, a; \theta)$  and then deriving the policy as  $\pi(s) := \arg \max_a Q(s, a; \theta)$ .
- ▶ **Gradients of the value function** can be computed from **gradients of the policy** via the **policy gradient theorem**:

$$\nabla_{\theta} V^{\pi}(s_0) \propto \sum_s \mu^{\pi}(s) \sum_a \nabla \pi(a | s; \theta) Q^{\pi}(s, a)$$

- ▶ Combined with observed values (Monte Carlo approach), gradient ascent for the (implicit) value function leads to a simple update rule for the policy parameters  $\theta$  called **REINFORCE**:

$$\theta := \theta + \alpha v_t \nabla \log \pi(a_t | s_t)$$

- ▶ any **baseline** for the observed value can be used and often accelerates convergence
  - ▶ esp. using a current estimate for the value function as baseline.

$$\theta := \theta + \alpha (v_t - \hat{V}(s_t)) \nabla \log \pi(a_t | s_t)$$

## Summary (2/2)

- ▶ Instead of using observed values (MC approach), one also can combine policy gradients with temporal differences (called **actor critic methods**):

$$\theta := \theta + \alpha(r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)) \nabla \log \pi(a_t | s_t)$$

- ▶ then besides the policy model (actor), a second model for the value function (critic) is required.

# Further Readings

- ▶ Reinforcement Learning:
  - ▶ Olivier Sigaud, Frederick Garcia (2010): *Reinforcement Learning*, ch. 1 in Sigaud and Buffet [2010].
  - ▶ Sutton and Barto [2018]

# References

Olivier Sigaud and Olivier Buffet, editors. *Markov Decision Processes in Artificial Intelligence*. Wiley, 2010.

David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *ICML*, 2014.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning*. MIT Press, 2nd edition edition, 2018.