

Reinforcement Learning

Lecture 11: Recap Session ¹

Lecturer: Prof. Dr. Mathias Niepert

Institute for Parallel and Distributed Systems
Machine Learning and Simulation Lab



University of Stuttgart
Germany

imprs-is

July 13, 2023

¹Many slides adapted from R. Sutton's course, D. Silver's course as well as previous RL courses given at U. of Stuttgart by J. Mainprice, D. Hennes, M. Toussaint, H. Ngo, and V. Ngo.

Outline

1. Evaluation
2. Study Suggestions
3. Value function approximation
4. Prediction
5. Control
6. Policy-based Reinforcement Learning
7. Policy Gradient Methods

Evaluation



<https://evasysw.uni-stuttgart.de/evasys/online.php?p=HP7FT>



<https://evasysw.uni-stuttgart.de/evasys/online.php?p=QSEE5>

Study Suggestions

What you need to study

The exam will cover topics in lectures 1-9

1. Multi-armed bandits; chapters: 2.1, 2.2, and 2.4 (+ softmax policy)
2. MDPs (definition, values function, etc.); chapters: 3.1-3.6
3. Policy improvement with dynamic programming; chapters 4.1-4.4
4. Monte-Carlo methods; chapter 5.1-5.7
5. Temporal difference methods; 6.1, 6.2, and 6.4-6.6
6. Planning and Learning; chapters: 8.1-8.5
7. Function approximation; chapter: 9.1-9.4, 9.5.4
8. n-step bootstrapping (no eligibility traces!); chapters: 7.1-7.3
9. Policy gradient methods; chapters: 13.1-13.6

General Remarks

- ▶ There will be no need to memorize or write pseudo-code (but pseudo-code typically helps understanding a method)
- ▶ You should memorize the core formulas corresponding to the various algorithms such as Bellman equation, MC, TD, etc.
- ▶ You should memorize backup diagrams of these formulas
- ▶ You may only use pen and scratch paper – no other materials (no textbooks, script, calculators, or mobiles) are allowed

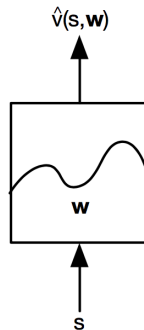
Value function approximation

Idea of value function approximation

- ▶ Parameterized functional form, with weights $\mathbf{w} \in \mathbb{R}^d$:

$$\hat{v}_{\pi}(s, \mathbf{w}) \approx v_{\pi}(s)$$

- ▶ Generally, much less weights than states $d \ll |\mathcal{S}|$
 - ▶ obvious for continuous state spaces
 - ▶ changing single weight, changes value estimate of many states
 - ▶ when one state is updated, change *generalizes* to many states
- ▶ Update \mathbf{w} with MC or TD learning



Function approximators for RL

- ▶ Differentiable function approximators, e.g.:
 - ▶ Linear combination of features
 - ▶ Neural networks
- ▶ RL specific problems:
 - ▶ non-stationary
 - ▶ non-iid data
 - ▶ bootstrapping
 - ▶ delayed targets

Random variables are *independent and identically distributed* (iid) if they each have the same probability distribution and are mutually independent

Stochastic Gradient Descent (SGD)

- ▶ Approximate value function $\hat{v}(s, \mathbf{w})$
 - ▶ differentiable for all $s \in \mathcal{S}$
- ▶ Weight vector $\mathbf{w} = (w_1, w_2, \dots, w_d)^\top$
 - ▶ \mathbf{w}_t weight vector at time $t = 0, 1, 2, \dots$
- ▶ Gradient of $f(\mathbf{w})$: $\nabla f(\mathbf{w}) = \left(\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right)^\top$
- ▶ Do gradient descent by sampling additive parts of the full gradient (i.e., each state consecutively)
- ▶ We can compute our update over smaller sets of inputs

Stochastic Gradient Descent (SGD) (part 2)

- ▶ When we **approximate** the gradient

- ▶ For example

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \mathcal{L}_n(\mathbf{w})$$

where \mathbf{w} are weights.

- ▶ In machine learning

$$\mathcal{L}(\mathbf{w}) = - \sum_{n=1}^N \log p(y_n \mid \mathbf{x}_n, \mathbf{w})$$

where $\mathbf{x}_n \in \mathbb{R}^D$ are training *inputs*
and y_n are the training *targets*

- ▶ The corresponding update

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha_t \sum_{n=1}^N (\nabla \mathcal{L}_n)(\mathbf{w}_t)$$

- ▶ Often the gradient is too difficult to compute (CPU/GPU expensive)

- ▶ **Mini-batch**: random subset

- a Large: accurate but costly
- b Small: noisy but cheap

Stochastic Gradient Descent (SGD) (part 3)

- ▶ *Mean squared Value Error*: $\mathcal{L}(\mathbf{w}) = \sum_{s \in S} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$
- ▶ Adjust \mathbf{w} to reduce the error on sample $S_t \mapsto v_\pi(S_t)$:

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{1}{2} \alpha_t (\nabla \mathcal{L}_t)(\mathbf{w}_t) \\ &= \mathbf{w}_t - \frac{1}{2} \alpha_t \nabla \underbrace{[v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2}_{\text{squared sample error}} \\ &= \mathbf{w}_t + \alpha_t [v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w})\end{aligned}$$

- ▶ α_t is a step size parameter
- ▶ Why not use $\alpha = 1$, thus eliminating full error on sample?

Linear methods

- ▶ Special case where $\hat{v}(\cdot, \mathbf{w})$ is *linear* in the weights
- ▶ **Feature vector** $\mathbf{x}(s)$ represents state s :

$$\mathbf{x}(s) = [x_1(s), \quad x_2(s), \quad \dots, \quad x_d(s)]^\top$$

- ▶ Each component of \mathbf{x} is a feature, examples:
 - ▶ distance of robot to landmarks
 - ▶ piece on a specific location on a chess board
- ▶ Value function is represented as a linear combination of features $\mathbf{x}(s)$:

$$\hat{v}(s, \mathbf{w}) = \mathbf{w}^\top \mathbf{x}(s) = \sum_{i=1}^d w_i x_i(s)$$

- ▶ Gradient is simply $\nabla \hat{v}(s, \mathbf{w}) = \mathbf{x}(s)$

Prediction

Prediction with function approximation

- ▶ We assumed the true value function $v_\pi(S_t)$ is known
- ▶ Substitute target U_t for $v_\pi(s)$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha[U_t - \hat{v}(S_t, \mathbf{w}_t)]\nabla\hat{v}(S_t, \mathbf{w}_t)$$

- ▶ U_t might be a noisy or bootstrapped approximation of the true value
- ▶ **Monte Carlo:** $U_t = G_t$
- ▶ **TD(0):** $U_t = R_{t+1} + \gamma\hat{v}(S_{t+1}, \mathbf{w}_t)$
- ▶ **TD(λ):** $U_t = G_t^\lambda$

Monte–Carlo with function approximation

- ▶ Target is unbiased by definition:

$$\mathbb{E}[U_t | S_t = s] = \mathbb{E}[G_t | S_t = s] = v_\pi(S_t)$$

- ▶ Training data:

$$\mathcal{D} = \{(S_1, G_1), (S_2, G_2), \dots, (S_{T-1}, G_{T-1}), (S_T, 0)\}$$

- ▶ Using SGD, \mathbf{w} is guaranteed to converge to a *local optimum*
- ▶ MC prediction exhibits local convergence with linear and non-linear function approximation
- ▶ SGD update for sample $S_t \mapsto G_t$:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha[G_t - \hat{v}(S_t, \mathbf{w})]\nabla \hat{v}(S_t, \mathbf{w})$$

Gradient Monte Carlo Algorithm

Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated

Input: a differentiable function $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameter: step size $\alpha > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop forever (for each episode):

 Generate an episode $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$ using π

 Loop for each step of episode, $t = 0, 1, \dots, T - 1$:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$

TD with function approximation

- ▶ TD-target $U_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$ is *biased* sample of the true value $v_\pi(S_t)$
- ▶ Training data:

$$\mathcal{D} = \{(S_1, R_2 + \gamma \hat{v}(S_2, \mathbf{w})), (S_2, R_3 + \gamma \hat{v}(S_3, \mathbf{w})), \dots, (S_{T-1}, R_T)\}$$

- ▶ SGD update for sample $S_t \mapsto G_t$:

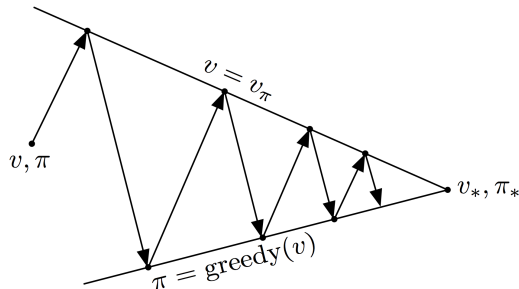
$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$$

- ▶ Linear TD(0):

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[\underbrace{R_{t+1} + \gamma \mathbf{w}^T \mathbf{x}(S_{t+1})}_{U_t: \text{TD-Target}} - \underbrace{\mathbf{w}^T \mathbf{x}(S_t)}_{\hat{v}: \text{value function}} \right] \mathbf{x}(S_t)$$

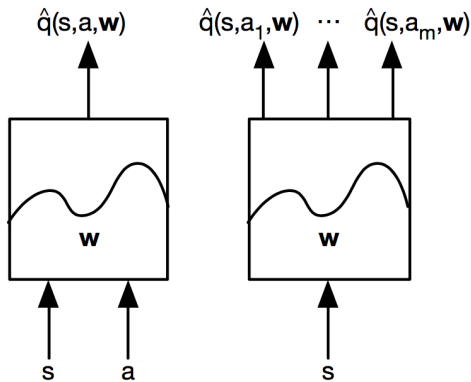
Control

Control with function approximation



- Control via generalized policy iteration (GPI):
 - *policy evaluation*: **approximate** policy evaluation: $\hat{q}(\cdot, \cdot, \mathbf{w}) \approx q_\pi$
 - *policy improvement*: ϵ -greedy policy improvement

Types of action–value function approximation



- ▶ Action as input: $\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$
- ▶ Multiple action–value outputs: $\hat{q}_a(s, \mathbf{w}) \approx q_\pi(s, a)$

Action-value function approximation

- ▶ Approximate action-value function $\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$
- ▶ Linear case:

$$\hat{q}(s, a, \mathbf{w}) = \mathbf{w}^T \mathbf{x}(s, a) = \sum_{i=1}^d w_i x_i(s, a)$$

$$\nabla \hat{q}(s, a, \mathbf{w}) = \mathbf{x}(s, a)$$

- ▶ Minimize squared error on samples $S_t, A_t \mapsto q_\pi$: $[q_\pi - \hat{q}(S_t, A_t, \mathbf{w})]^2$
- ▶ Use SGD to find local minimum:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{1}{2} \alpha \nabla [q_\pi(S_t, A_t) - \hat{q}(S_t, A_t, \mathbf{w}_t)]^2 \\ &= \mathbf{w}_t + \alpha [q_\pi(S_t, A_t) - \hat{q}(S_t, A_t, \mathbf{w}_t)] \nabla \hat{q}(S_t, A_t, \mathbf{w}) \end{aligned}$$

Control with function approximation

- ▶ Again, we must substitute target U_t for true action-value $q_\pi(s, a)$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha[U_t - \hat{q}(S_t, A_t, \mathbf{w}_t)]\nabla\hat{q}(S_t, A_t, \mathbf{w})$$

- ▶ **Monte Carlo:** $U_t = G_t$
- ▶ **One-step Sarsa:** $U_t = R_{t+1} + \gamma\hat{q}(S_{t+1}, A_{t+1}, \mathbf{w})$

Policy-based Reinforcement Learning

Policy-based reinforcement learning

- ▶ So far we approximated the action-value function and generated a policy from it
- ▶ Approximation of the action-value function

$$\hat{q}(s, a, \mathbf{w}) \approx q_{\pi}(s, a)$$

- ▶ Generation of policy by, e.g., ϵ -greedy

$$\hat{q}(s, a, \mathbf{w}) \xrightarrow{\epsilon\text{-greedy}} \pi$$

- ▶ Now we directly *parameterize* the policy π

Policy optimization

- Policy optimization:

$$\pi_* = \pi(a \mid s, \theta_*)$$

with

$$\theta_* = \arg \max_{\theta} J(\theta)$$

where J is some *performance measure*

- Discounted return: $G_0 = \sum_{k=0}^{T-1} \gamma^k R_{k+1}$

$$\pi_* = \arg \max_{\pi} \mathbb{E}_{\pi}[G_0] = \arg \max_{\pi} \mathbb{E}_{\pi}[v_{\pi}(s_0) \mid S_0 = s_0]$$

- Undiscounted return: $G_0 = \sum_{k=0}^{T-1} R_{k+1}$, i.e., $\gamma = 1$

$$\pi_* = \arg \max_{\pi} \mathbb{E}_{\pi}[R_0 + R_1 + \dots + R_{T-1} \mid S_0 = s_0]$$

Policy optimization

- In continuing environments, we could use the **average** state-value:

$$\sum_s \mu_\pi(s) \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) r(s, a, s')$$

where μ_π is the steady-state distribution under π . This is useful for the function approximation where states are not well defined (i.e., only seen through their features).

Parameterized policies

- ▶ Policies parameterized by parameter $\theta \in \mathbb{R}^d$:
 - ▶ deterministic: $a = \pi(s, \theta)$
 - ▶ stochastic: $a \sim \pi(a \mid s, \theta) = \Pr\{A_t = a \mid S_t = s, \theta\}$
- ▶ Objective becomes $\max_{\theta} J(\theta) = \max_{\theta} \mathbb{E}_{\pi_{\theta}}[G_0]$
- ▶ We can parameterize π in any way, as long as it is *differentiable* wrt to θ
- ▶ In general we require $\pi(a \mid s, \theta) \in [0, 1]$ for all s, a
- ▶ If the action space is discrete (and not too large): **softmax policy**

$$\pi(a \mid s, \theta) = \frac{e^{h(s,a,\theta)}}{\sum_b e^{h(s,b,\theta)}}$$

where $h(\cdot)$ is the *action preference* function

- ▶ Can this approach the deterministic policy?

Policy gradient methods

► Problem:

$$\pi_* = \pi(a \mid s, \theta_*)$$

with

$$\theta_* = \arg \max_{\theta} J(\theta) = \arg \max_{\theta} \mathbb{E}_{\pi_{\theta}}[G_0]$$

Intuition: collect a bunch of trajectories, and ...

1. Make the good trajectories more probable
2. Make the good actions more probable
3. Push the policy towards generating good actions

► Policy gradient methods:

$$\theta \leftarrow \theta + \alpha \widehat{\nabla J(\theta)}$$

- where $\nabla J(\theta)$ is the *policy gradient*:

$$\nabla J(\theta) = \left(\frac{\partial J(\theta)}{\partial \theta_0}, \dots, \frac{\partial J(\theta)}{\partial \theta_d} \right)^T$$

Score function

- ▶ We now compute the policy gradient analytically
- ▶ Assume policy π_{θ} is differentiable whenever it is non-zero
- ▶ and we know the gradient $\nabla_{\theta} \pi_{\theta}(a | s)$
- ▶ We have the following useful identity

$$\begin{aligned}\nabla_{\theta} \pi_{\theta}(a | s) &= \pi_{\theta}(a | s) \frac{\nabla_{\theta} \pi_{\theta}(a | s)}{\pi_{\theta}(a | s)} \\ &= \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s)\end{aligned}$$

- ▶ The **score function** is $\nabla_{\theta} \log \pi_{\theta}(a | s)$

Score: Softmax Policy

- ▶ Weight actions using linear combination of features $h(s, a) = \phi(s, a)^\top \theta$
- ▶ Probability of action is proportional to exponentiated weight

$$\pi_\theta \sim \exp\left(\phi(s, a)^\top \theta\right)$$

- ▶ The score function is

$$\nabla_\theta \log \pi_\theta(a \mid s) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, \cdot)]$$

Score: Gaussian Policy

- ▶ In continuous action spaces, a Gaussian policy is natural
- ▶ Mean is a linear combination of state features $\mu(s) = \phi(s)^\top \theta$
- ▶ Variance may be fixed σ^2 , or can also be parametrized
- ▶ Policy is Gaussian, $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- ▶ The score function is

$$\nabla_{\theta} \log \pi_{\theta}(a | s) = \frac{(a - \mu(s)) \phi(s)}{\sigma^2}$$

One-Step MDPs

- ▶ Consider a simple class of one-step MDPs
 - ▶ Starting in state $s \sim \mu(s)$
 - ▶ Terminating after one time-step (taking action a) with reward $r = R_{s,a}$
- ▶ Use score function to compute the policy gradient

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\pi_{\theta}}[r] \\ &= \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(a \mid s) R_{s,a} \\ \nabla_{\theta} J(\theta) &= \sum_{s \in \mathcal{S}} \mu(s) \sum_{a \in \mathcal{A}} \pi_{\theta}(a \mid s) \nabla_{\theta} \log \pi_{\theta}(a \mid s) R_{s,a} \\ &= \mathbb{E}_{\pi_{\theta}}[r \nabla_{\theta} \log \pi_{\theta}(a \mid s)] \end{aligned}$$

Policy Gradient Theorem

- ▶ The policy gradient theorem generalises the previous derivation to multi-step MDPs
- ▶ Replaces instantaneous reward r with long-term value $q_\pi(s, a)$
- ▶ Policy gradient theorem applies to start state objective, average reward and average value objective

Theorem (Policy-gradient)

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} [q_{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \log \pi(a \mid s, \boldsymbol{\theta})]$$

Policy Gradient Methods

Monte–Carlo policy gradient (REINFORCE)

- Update parameters θ by stochastic gradient ascent
 - using the policy gradient theorem
 - using return G_t as an unbiased sample of $q_\pi(S_t, A_t)$

Initialize a differentiable policy parameterization $\pi(a|s, \theta)$

Initialize $\theta \in \mathbb{R}^{d'}$

for all episodes **do**

 Generate an episode $\tau = (S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T)$
 following $\pi(\cdot | \cdot, \theta)$

for $t = 0, 1, \dots, T - 1$ **do**

$G_t \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$ // return at time t

$\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \log \pi(A_t | S_t, \theta)$ // update rule

end for

end for

REINFORCE with baseline

- Monte–Carlo policy gradient still suffers from high variance

Theorem (Policy-gradient with baseline)

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left[(q_{\pi}(s, a) - b(s)) \nabla_{\boldsymbol{\theta}} \log \pi(a \mid s, \boldsymbol{\theta}) \right]$$

Update rule: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t (G_t - b(S_t)) \nabla_{\boldsymbol{\theta}} \log \pi(A_t \mid S_t, \boldsymbol{\theta})$

- Still unbiased with lower variance!

REINFORCE with baseline

- ▶ We use an approximation of the value function $\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$ as a *baseline*
- ▶ Policy parameters: $\boldsymbol{\theta}$
- ▶ Value estimator parameters: \mathbf{w}
- ▶ Update rule:

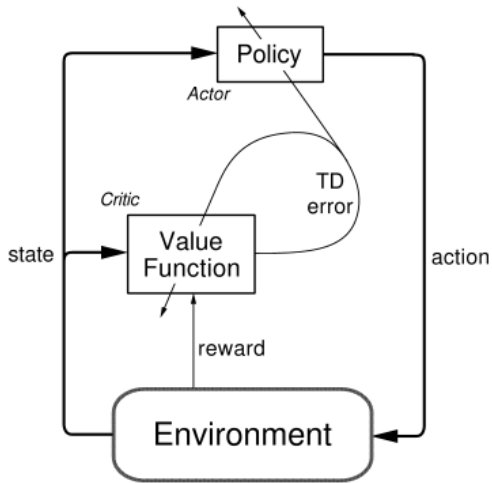
$$\begin{aligned}\delta &\leftarrow G - \hat{v}(S_t, \mathbf{w}) \\ \mathbf{w} &\leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \gamma^t \delta \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}) \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} \gamma^t \delta \nabla_{\boldsymbol{\theta}} \log \pi(A_t | S_t, \boldsymbol{\theta})\end{aligned}$$

- ▶ Interpretation:
 - ▶ \uparrow log-prob of action A_t proportionally to how much G_t is better than expected
 - ▶ baseline accounts for and removes the effect of past actions

Estimating the action–value function

- ▶ We are solving the prediction problem: policy evaluation
- ▶ How good is policy π_{θ} with current parameters θ ?
- ▶ Familiar toolset for *fitting* the baseline:
 - ▶ Monte–Carlo policy evaluation
 - ▶ TD-learning
 - ▶ TD(λ)
 - ▶ LSPI

Actor-critic concept



Actor-critic vs. baseline

- ▶ Actor-critic methods use the value function as a baseline for policy gradients
- ▶ Delivers trade off between *variance reduction* of policy gradients with *bias introduction* from value function methods
- ▶ **Critic:** updates value-function parameters w
- ▶ **Actor:** updates policy parameters θ using critic
- ▶ REINFORCE with baseline uses value-function as *baseline* not as a *critic*
 - ▶ not used for *bootstrapping*
- ▶ One-step actor-critic update:

$$\begin{aligned}\delta &\leftarrow R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w) \\ w &\leftarrow w + \alpha^w \gamma^t \delta \nabla_w \hat{v}(S_t, w) \\ \theta &\leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla_\theta \log \pi(A_t | S_t, \theta)\end{aligned}$$

Bias in Actor–Critic algorithms

- ▶ Approximating (bootstrapping) the policy gradient introduces bias
- ▶ Biased policy gradient might not find the right solution
- ▶ But reduces variance and makes learning substantially more efficient
- ▶ *Compatible function approximation* avoids this problem: