# Lecture 11: Off-policy and multi-step learning

Hado van Hasselt

UCL, 2021

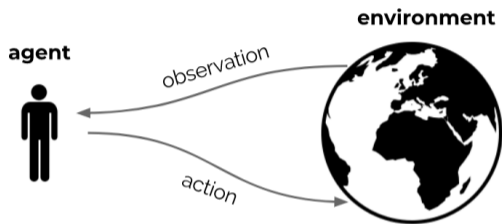# Background

Sutton & Barto 2018, Chapter 5, 7, 11

# Recap

# Recap



- Reinforcement learning is the science of learning to make decisions
- Agents can learn a **policy**, **value function** and/or a **model**
- The general problem involves taking into account **time** and **consequences**
- Decisions affect the **reward**, the **agent state**, and **environment state**

# High level

- Previous lectures:
  - Model-free prediction & control
  - Multi-step updates (and eligibility traces)
  - Understanding dynamic programming operators
  - Predictions with function approximation
  - Model-based algorithms
  - Policy gradients and actor-critic algorithms
- This lecture:
  - **Off-policy learning**, especially when combined with **multi-step** updates and **function approximation**

# Motivation

# Why learn off-policy?

Why learn off-policy?

- Off-policy learning is important to learn about **hypothetical**, **counterfactual** events (i.e, "what if" question)

- Use cases include

  - learning about the greedy policy

  - learning about (many) other policies

  - learning from observed data (e.g., stored logs / other agents)

  - learning from past policies

- This is also important to correct for **mismatch in data distributions**, for instance for policy gradients

# One-step off-policy

With action values, **one-step** off-policy learning seems relatively straightforward:

$$q(S_t, A_t) \leftarrow q(S_t, A_t) + \alpha_t(R_{t+1} + \sum_a \pi(a|S_{t+1})q(S_{t+1}, a) - q(S_t, A_t))$$

For instance

▶ Q-learning: let $\pi$ be greedy $\implies$ $\sum_a \pi_{sa}q_{sa} = \max_a q_{sa}$

▶ Expected Sarsa: let $\pi$ be the current behaviour policy

▶ Sarsa: let $\pi$ put all probability mass on the action the behaviour picked

# Multi-step off-policy

For **multi-step updates**, we can use **importance-sampling corrections**
E.g., for a Monte Carlo return on a trajectory $\tau_t = \{S_t, A_t, R_{t+1}, \ldots, S_T\}$

$$\hat{G}_t \equiv \frac{p(\tau_t | \pi)}{p(\tau_t | \mu)} G_t = \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} \cdots \frac{\pi(A_T | S_T)}{\mu(A_T | S_T)} G_t,$$
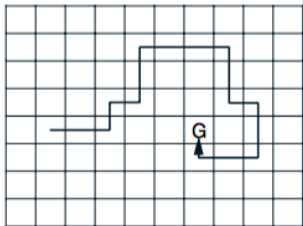
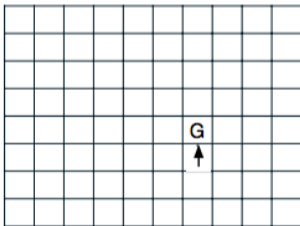then $\mathbb{E}[\hat{G}_t \mid \mu] = \mathbb{E}[G_t \mid \pi]$

# Multi-step off-policy

▶ We know multi-step updates often more efficiency propagate information

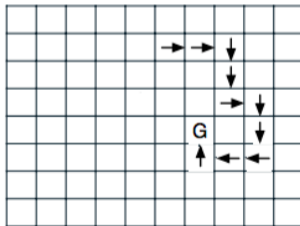▶ Neither full Monte Carlo, nor one-step bootstrapping, is typically the best trade-off



Path taken

Action values increased by one-step Sarsa

Action values increased by 10-step Sarsa

# Off-policy corrections for policy gradients

- Recall: for policy gradient methods we want to sample/estimate

$$\mathbb{E}[q_\pi(S_t, A_t)\nabla \log \pi(A_t|S_t)].$$

- On-policy can sample multi-step returns $G_t$ such that $\mathbb{E}[G_t \mid \pi] \approx q_\pi(s, a)$
- But what if the behaviour is $\mu \neq \pi$?
- $\implies$ we might not be following a gradient direction

# Issues in off-policy learning

# Issues with off-policy learning

The following issues (especially) arise when learning off-policy

- ▶ High variance (especially when using multi-step updates)
- ▶ Divergent and inefficient learning (especially when using one-step updates)

We will discuss both in this lecture

# Issues in off-policy learning: Variance

# Variance of importance sampling corrections

- A big issue in using importance-sampling corrections is **high variance**
- First, consider a one-step reward
- Verify the expectation, for a given state $s$:

$$\mathbb{E}\left[\frac{\pi(A_t|s)}{\mu(A_t|s)}R_{t+1} \mid A_t \sim \mu\right] = \sum_a \mu(a|s)\frac{\pi(a|s)}{\mu(a|s)}r(s,a)$$
$$= \sum_a \pi(a|s)r(s,a)$$
$$= \mathbb{E}[R_{t+1} \mid A_t \sim \pi]$$

- But typically the variance will be larger, sometimes greatly so
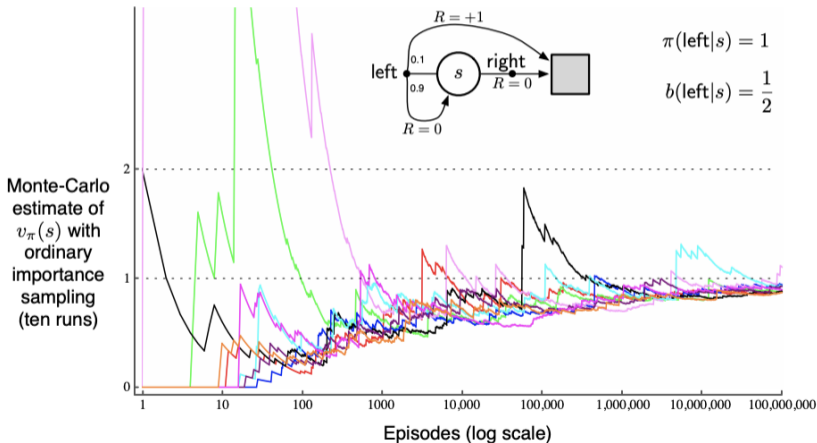
# Variance example

# Variance of importance sampling corrections

- In some cases the variance of an importance-weighting return can even be **infinite** (see: Sutton & Barto, Example 5.5)

# Mitigating variance

- There are multiple ways to reduce variance
- We will discuss three:
    - Per-decision importance weighting
    - Control variates
    - (Adaptive) bootstrapping

# Reducing variance:
# Per-decision importance weighting

# Mitigating variance: with per-decision importance weighting

Consider some state $s$. For any random $X$ that does not correlate with (random) action $A$ we have

$$\mathbb{E}[X \mid \pi] = \mathbb{E}\left[\frac{\pi(A|s)}{\mu(A|s)}X \mid \mu\right] = \mathbb{E}[X \mid \mu]$$

Intuition: the expectation does not depend on the policy, so we don't need to correct

# Mitigating variance: with per-decision importance weighting

Proof:

$$\mathbb{E}\left[\frac{\pi(A|s)}{\mu(A|s)}X \mid \mu\right]$$

$$= \mathbb{E}[X \mid \mu]\mathbb{E}\left[\frac{\pi(A|s)}{\mu(A|s)} \mid \mu\right] \qquad \text{(Because } X \text{ and } \frac{\pi}{\mu} \text{ are uncorrelated)}$$

$$= \mathbb{E}[X \mid \mu]\sum_a \mu(a|s)\frac{\pi(a|s)}{\mu(a|s)}$$

$$= \mathbb{E}[X \mid \mu]\sum_a \pi(a|s)$$

$$= \mathbb{E}[X \mid \mu] \qquad \text{(Because } \sum_a \pi(a|s) = 1)$$

Similarly, in general, we have $\mathbb{E}[\frac{\pi(A|s)}{\mu(A|s)} \mid \mu] = 1$

# Notation

Shorthand notations:

$$\rho_t \equiv \frac{\pi(A_t \mid S_t)}{\mu(A_t \mid S_t)} \qquad\qquad \rho_{t:t+n} \equiv \prod_{k=t}^{t+n} \rho_k = \prod_{k=t}^{t+n} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$$

Then the reweighted MC return from state $S_t$ terminating at time $T$ can be written as

$$\overbrace{\left( \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} \right)}^{= \ \rho_{t:T-1}} \overbrace{\left( \sum_{k=t}^{T-1} \gamma^{k-t} R_{k+1} \right)}^{= \ G_t} = \rho_{t:T-1} G_t = \sum_{k=t}^{T-1} \rho_{t:T-1} \gamma^{k-t} R_{k+1}$$

We can interpret the importance-weight $\rho_{t:T-1}$ as applying to each reward

# Mitigating variance: with per-decision importance weighting

$$\rho_{t:T-1} G_t = \sum_{k=t}^{T-1} \rho_{t:T-1} \gamma^{k-t} R_{k+1}$$

**Earlier** rewards cannot depend on **later** actions. This means:

$$\mathbb{E}[\rho_{t:T-1} G_t \mid \mu] = \mathbb{E}[\sum_{k=t}^{T-1} \rho_{t:T-1} \gamma^{k-t} R_{k+1} \mid \mu]$$

$$= \mathbb{E}[\sum_{k=t}^{T-1} \rho_{t:k} \gamma^{k-t} R_{k+1} \mid \mu]$$

Recursive definition of the latter:

$$G_t^\rho = \rho_t(R_{t+1} + \gamma G_{t+1}^\rho)$$

# Mitigating variance: with per-decision importance weighting

▶ Per-decision importance-weighted return

$$G_t^\rho = \rho_t(R_{t+1} + \gamma G_{t+1}^\rho)$$

We can use this to learn $v_\pi$ from data generated under $\mu \neq \pi$

▶ To learn **action values** $q_\pi$, we can use

$$G_t^\rho = R_{t+1} + \gamma \rho_{t+1} G_{t+1}^\rho$$

▶ How and why are these different?

# Reducing variance: Control variates

# Example: control variates

# Control variates for multi-step returns

The idea of control variates can be extended to multi-step returns

► First, recall

$$
\begin{aligned}
\delta_t^\lambda &\equiv G_t^\lambda - v(S_t) \\
&= R_{t+1} + \gamma((1 - \lambda)v(S_{t+1}) + \gamma\lambda G_{t+1}^\lambda) - v(S_t) \\
&= \underbrace{R_{t+1} + \gamma v(S_{t+1}) - v(S_t)}_{= \delta_t} + \gamma\lambda \underbrace{(G_{t+1}^\lambda - v(S_{t+1}))}_{= \delta_{t+1}^\lambda} \\
&= \delta_t + \gamma\lambda\delta_{t+1}^\lambda
\end{aligned}
$$

# Control variates for multi-step returns

The idea of control variates can be extended to multi-step returns

- ▶ Now, lets add per-decision importance weights

$$\delta_t^\lambda = \delta_t + \gamma\lambda\delta_{t+1}^\lambda$$
$$\delta_t^{\rho\lambda} = \rho_t(\delta_t + \gamma\lambda\delta_{t+1}^{\rho\lambda})$$

- ▶ By design this includes the $(1 - \rho_t)v(S_t)$ control variate terms
- ▶ Sometimes called 'error weighting' (to contrast to 'reward weighting')

# Control variates for multi-step returns

$$\delta_t^{\rho\lambda} = \rho_t(\delta_t + \gamma\lambda\delta_{t+1}^{\rho\lambda})$$

▶ One can show that

$$\mathbb{E}[\delta_t^{\rho\lambda} \mid \mu] = \mathbb{E}[G_t^{\rho\lambda} - v(S_t) \mid \mu]$$

where

$$G_t^{\rho\lambda} = \rho_t\left(R_{t+1} + \gamma\left((1 - \lambda)v(S_{t+1}) + \lambda G_{t+1}^{\rho\lambda}\right)\right)$$

is the per-decision importance-weighted $\lambda$-return.

▶ But $\delta_t^{\rho\lambda}$ can have lower variance than $G_t^{\rho\lambda} - v(S_t)$

# Reducing variance: Adaptive Bootstrapping

# Reducing variance: bootstrapping

- ► For our last technique, we consider bootstrapping
- ► This amounts to picking $\lambda < 1$ when using either $\delta_t^{\rho\lambda}$ or $G_t^{\rho\lambda}$
- ► Note that to learn **action values**, we can use

$$G_t^{\rho\lambda} = R_{t+1} + \gamma\left((1-\lambda)\sum_a \pi(a \mid S_{t+1})q(S_{t+1}, a) + \rho_{t+1}\lambda G_{t+1}^{\rho\lambda}\right)$$

Then, if $\lambda = 0$, we get

$$G_t = R_{t+1} + \gamma\sum_a \pi(a \mid S_{t+1})q(S_{t+1}, a)$$

$\implies$ no more importance weighted $\implies$ low variance

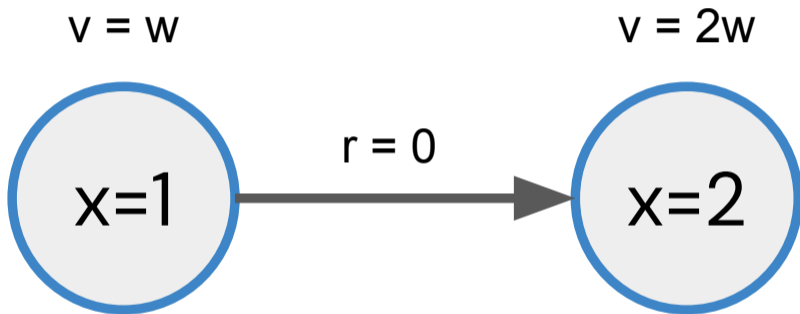- ► However, bootstrapping too much may open us to the **deadly triad**!

# Recap: Deadly triad

▶ Recall, the deadly triad refers to the possibility of divergence when we combine
  ▶ **Bootstrapping**
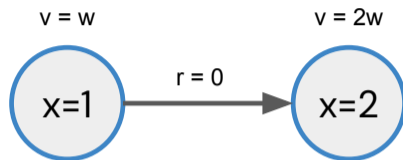  ▶ **Function approximation**
  ▶ **Off-policy learning**

# Recap: Deadly triad



What if we use TD only on this transition?

# Recap: Deadly triad



v = w          v = 2w

x=1    r = 0    x=2

$$w_{t+1} = w_t + \alpha_t(r + \gamma v(s') - v(s))\nabla v(s)$$
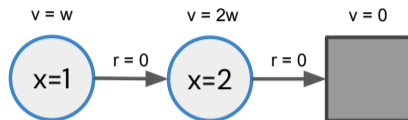$$= w_t + \alpha_t(2\gamma - 1)w_t$$

Suppose $\gamma > \frac{1}{2}$. Then,
When $w_t > 0$ and , then $w_{t+1} > w_t$
When $w_t < 0$ and , then $w_{t+1} < w_t \implies w_t$ diverges to $+\infty$ of $-\infty$

# Recap: Deadly triad



- What if we use multi-step returns?
- Still consider only updating the left-most state

$$\Delta w = \alpha(r + \gamma(G_t^\lambda - v(s)))$$
$$= \alpha(2\gamma(1 - \lambda) - 1)w$$

- The multiplier is negative when $2\gamma(1 - \lambda) < 1 \implies \lambda > 1 - \frac{1}{2\gamma}$
- E.g., when $\gamma = 0.9$, then we need $\lambda > 4/9 \approx 0.45$
- Conclusion: if we do not bootstrap too much, we can learn better

# Reducing variance: **adaptive** bootstrapping

- We don't want to bootstrap too much $\implies$ **deadly triad**
- We don't want to bootstrap too little $\implies$ **high variance**
- Can we adaptively bootstrap 'just enough'?
- Idea: bootstrap **adaptively** only in as much as you go off-policy

# Reducing variance: adaptive bootstrapping

- Recall $\delta_t^{\rho\lambda} = \rho_t(\delta_t + \gamma\lambda\delta_{t+1}^{\rho\lambda})$
- Let's add an initial bootstrap parameter, and make these time-dependent

$$\delta_t^{\rho\lambda} = \lambda_t\rho_t(\delta_t + \gamma\delta_{t+1}^{\rho\lambda})$$

  (If $\lambda_t = 1$, we obtain the previous version)

- We can pick $\lambda_t$ **separately on each time step**
- Idea: pick it such that, for all $t$, $\lambda_t\rho_t \leq 1$:

$$\lambda_t = \min(1, 1/\rho_t)$$

- Intuition: when we are too off-policy ($\rho$ is far from one) truncate the sum of errors
- This is the same as bootstrapping there

# Reducing variance: adaptive bootstrapping

$$\lambda_t = \min(1, 1/\rho_t)$$

▶ This is known as **ABTD** (Mahmood et al. 2017) or **v-trace** (Espeholt et al. 2018)

▶ We are free to choose different ways to bootstrap: in the tabular case all these methods will be updating towards some mixture of multi-step returns, and therefore converge

▶ In deep RL this really helps, especially for policy gradients
(Policy gradients do not like biased return estimates – we will get back to that)

▶ This is used a lot these days

# Reducing variance: tree backup

- Picking $\lambda_t = \min(1, 1/\rho_t)$ is not the only way to adaptively bootstrap

- One more option, consider the Bellman operator for action values

$$q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})q_\pi(S_{t+1}, a) \mid A_t = a, S_t = s]$$

- Note: the expectation does not depend on $\pi$, because we condition on the action $a$

- Idea: sample this, then replace only the action you selected:

$$G_t = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})q(S_{t+1}, a) + \boldsymbol{\gamma\pi(A_{t+1}|S_{t+1})G_{t+1}}$$

- We remove only the expectation $q(S_{t+1}, A_{t+1})$ of the action actually selected, and replace it with the return

- This is **unbiased**, and **low variance**! ($\pi(A_{t+1}|S_{t+1})$ plays a role similar to $\lambda$)

- It might bootstrap too early though — beware of deadly triads!

# Lecture End