

# Reinforcement Learning

## – Course notes –

December 14, 2019

# Contents

<b>1</b>	<b>Markov Decision Processes</b>	<b>2</b>
1.1	Decision rules, Policies	2
<b>2</b>	<b>Value functions, Dynamic programming</b>	<b>3</b>
2.1	The value function	3
2.1.1	The optimal value function	4
2.2	The $Q$ -function	5
2.2.1	The optimal $Q$ -function	5
2.3	Algorithms for dynamic programming	6
2.3.1	Value Iteration	6
2.3.2	Policy Iteration	7
<b>3</b>	<b>Approximate Reinforcement Learning</b>	<b>8</b>
3.1	Monte Carlo Policy Evaluation	8
3.1.1	Incremental MC	9
3.2	Temporal difference (TD) estimation	9
3.2.1	Using the Bellman equation: TD(0) estimation	9
3.2.2	Temporal difference TD( $\lambda$ )	10
3.3	Application of TD to policy improvement	10
3.3.1	SARSA	10
3.3.2	Q-learning	11
<b>4</b>	<b>Approximate solving of Markov Decision Processes with Policy Gradients</b>	<b>12</b>
4.1	Monte Carlo policy gradient: the REINFORCE algorithm	13
4.1.1	Variance reduction: temporal structure	13
4.1.2	Variance reduction with baselines	14
4.1.3	Parametric Bootstrapping of the baseline	15
4.2	Bootstrapped returns: Actor-Critic algorithms	16
4.2.1	Actor-Critic	16
4.2.2	Actor-Critic with baselines: Advantage Actor-Critic (A2C)	17
<b>5</b>	<b>The exploration-exploitation dilemma</b>	<b>19</b>
5.1	Multi-Armed Bandits	19

# Chapter 1

## Markov Decision Processes

**Definition 1 (Markov Decision process)** *A Markov decision process (MDP) is given by its state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , state dynamics, reward structure and discount factor. In a MDP, the state dynamics can be written as*

$$\mathbb{P}(s_{t+1} \mid s_t, a_t)$$

The **return** of a trajectory  $(s_0, a_0, r_0, s_1, \dots)$  at time  $t$  is defined as

$$G_t = \sum_k \gamma^k r_{t+k} \tag{1.1}$$

### 1.1 Decision rules, Policies

The idea of a *decision rule* is to choose which action (or actions, with a given preference) to take after observing part of an MDP's trajectory.

# Chapter 2

## Value functions, Dynamic programming

### 2.1 The value function

The value function is a staple from the literature on dynamic programming, whether it be for discrete or continuous problems (as in control theory). It measures just how good a control  $u$  – or, in our case, a policy  $\pi$  – is regarding the desired target of our problem.

**Definition 2 (Value function)** *The value function  $V^\pi: \mathcal{S} \rightarrow \mathbb{R}$  of a policy  $\pi$  is the expectation of the cumulative (discounted) future rewards starting from a point  $s_0 = s$*

$$V^\pi(s) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \mid s_0 = s \right] \quad (2.1)$$

where the trajectory  $\tau$  is generated under the policy  $\pi$ .

$T$  can be a random stopping time corresponding to the first time of arrival at a so-called terminal state. It can also be  $T = \infty$  for infinite horizon problems.

**Remark 1 (Notational abuse)** *Often, we will write the value function at a state  $s_t$ , where it is implied we are at the  $t$ -th step in a trajectory, as*

$$V(s_t) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}, a_{t'}) \mid s_t \right]$$

*This is coherent with the notion of cumulative (discounted) future rewards which defined the value function (2.1), but not with the notation. In the finite horizon setting, a more correct notation would be to write  $V(t, s_t)$  and make the dependence on the starting time  $t$  of the discounting explicit.*

This notion can be generalized to other cases, such as the case where the rewards are generated by the transitions  $(s_t, a_t, s_{t+1})$  rather than the (state, action) couple:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t, s_{t+1}) \mid s_0 = s \right]$$

or the stochastic reward case where  $r_t$  is distributed according to some law parameterized by the state and action.

Under a deterministic policy  $\pi: \mathcal{S} \rightarrow \mathcal{A}$  and associated decision rule  $d^\pi(s) = \pi(s)$ , the dynamic programming principle leads to a dynamic programming equation called the **Bellman equation**:

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s, \pi(s), s') V^\pi(s') \quad (2.2)$$

which is a fixed-point condition.

The Bellman equation can be used to *evaluate* a policy  $\pi$ , that is compute its associated value function  $V^\pi$ .

Its fixed-point structure can be reformulated in terms of an operator on a function space, called the **Bellman operator**:

$$\mathcal{T}^\pi v(s) := r(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s, \pi(s), s') v(s')$$

which can be shown to be contractant, ensuring the existence of a solution.

**Remark 2** *We have the following, possible generalizations (see Sutton's book [1, chap. 3, 4] for further details):*

- for stochastic policies  $\pi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$ , the sum in eq. (2.2) becomes

$$\sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi(s, a) p(s, a, s') V^\pi(s')$$

- in non-discrete state spaces, the sum can be replaced by an integral with respect to a measure  $p(s, \pi(s), ds')$
- if the rewards are given for transitions as  $r(s, a, s')$ , we introduce

$$r(s, a) = \sum_{s' \in \mathcal{S}} p(s, a, s') r(s, a, s'),$$

and the Bellman equation can be rewritten

$$V^\pi(s) = \sum_{s' \in \mathcal{S}} p(s, \pi(s), s') (r(s, \pi(s), s') + \gamma V^\pi(s'))$$

In a finite setting with small state and action spaces, it possible to solve the Bellman equation directly by Gaussian elimination on the corresponding matrices and vectors.

### 2.1.1 The optimal value function

Solving a Markov Decision Process involves finding an **optimal** policy  $\pi^*$  that will maximize the expected rewards in the long run when starting from a given state  $s_0$  (or distribution  $s_0 \sim p$ ).

**Definition 3 (Optimal policy and value function)** *Given a set of policies  $\Pi$ , the **optimal value function** satisfies*

$$V^* = \max_{\pi} V^\pi \quad (2.3)$$

*for every state  $s \in \mathcal{S}$ . An **optimal policy**  $\pi^*$  is one that satisfies the maximum.*

Strictly speaking, we are taking a maximal policy with respect to a partial ordering on policies that compares them by looking at their respective value functions:

$$\pi_1 \leq \pi_2 \iff [V^{\pi_1}(s) \leq V^{\pi_2}(s) \ \forall s \in \mathcal{S}]$$

**Proposition 1 (Optimal Bellman equation)** *The optimal value function  $V^*$  obeys a dynamic programming principle, the **optimal Bellman equation**.*

$$V^*(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \sum_{s' \in \mathcal{S}} p(s, a, s') V^*(s') \right\} \quad (2.4)$$

This is also a fixed-point condition, which can once again be expressed in terms of an operator called the **Bellman optimal operator**:

$$\mathcal{T}^* v(s) := \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a, s') v(s') \right\}.$$

Then, we have that the optimal value function satisfies the equation

$$\mathcal{T}^* V^* = V^*$$

## 2.2 The $Q$ -function

**Definition 4 (action-state value function)** *The action-state value function of a policy  $\pi$  is the function  $Q^\pi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is defined by*

$$Q^\pi(s, a) := \mathbb{E}_\pi \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right] \quad (2.5)$$

where the trajectory  $\tau$  is generated under the policy  $\pi$ . The horizon  $T$  of the problem can be finite or infinite ( $T$  can be a stopping time).

The action-state value function  $Q^\pi$  has an obvious link to the value function  $V^\pi$ . For any state  $s \in \mathcal{S}$ , it holds that

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s, \cdot)} [Q^\pi(s, a)] = \sum_a \pi(s, a) Q^\pi(s, a) \quad (2.6)$$

There's also a reverse equality:

$$Q^\pi(s, a) = \mathbb{E}_\pi [r(s, a) + V^\pi(s_{t+1}) \mid s_t = s] = r(s, a) + \sum_{s'} p(s, a, s') V^\pi(s') \quad (2.7)$$

The action-state value function has its own **Bellman equation**:

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s, a, s') Q^\pi(s', a) \quad (2.8)$$

### 2.2.1 The optimal $Q$ -function

**Definition 5 (Optimal action value function)** *Optimal policies as defined in section 2.1.1 also share the same **optimal action-state value function***

$$Q^* = \max_\pi Q^\pi \quad (2.9)$$

for all states  $s \in \mathcal{S}$  and actions  $a \in \mathcal{A}$ .

Here, the maximum is also understood with respect to a partial ordering on policies.

The optimal value  $Q^*$  follows its own variant of the **optimal Bellman equation**:

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a, s') \max_{a'} Q^*(s', a') \quad (2.10)$$

### Link between the optimal value functions

We also have a link between the optimal value functions. Using the Bellman optimality equation on  $Q$  (2.10), we get:

$$Q^*(s, a) = \mathbb{E}_\pi [r_t + \gamma V^*(s_{t+1}) \mid s_t = s, a_t = a]$$

Conversely, the optimal value function can be expressed by reformulating the Bellman optimality equation (2.4):

$$V^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a) \quad (2.11)$$

### Obtaining the optimal policy

Given the optimal action-state value function  $Q^*$ , you can define an optimal policy by taking

$$\pi^*(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a).$$

By the way, this shows we can always find a deterministic optimal policy. This is an **important** property of the optimal  $Q$ -function which we will use to design policy optimization algorithms.

## 2.3 Algorithms for dynamic programming

### 2.3.1 Value Iteration

Under suitable hypotheses, it can be shown that the Bellman operators  $\mathcal{T}^\pi$  and  $\mathcal{T}^*$  are contractions with respect to the infinity norm with Lipschitz constant  $\gamma$ .

This means we can approximate the optimal value function by iterating the optimal Bellman operator  $\mathcal{T}^*$ . This leads to the following *value iteration* algorithm:

---

#### Algorithm 2.1: Value iteration

---

**Input:** Rewards  $r(s, a)$ , transitions  $p(s, a, s')$ , initial value proposal  $V_0$ , number of iterations  $K$

1 **foreach**  $k = 1, \dots, K$  **do**

2    $V_k \leftarrow \mathcal{T}^* V_{k-1}$ ;

3 **foreach**  $s \in \mathcal{S}$  **do**

4    $\pi_K(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \{r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a, s') V_K(s')\}$ ;

5 **return** Policy  $\pi_K$ , value  $V_K$ ;

---

**Remark 3** Each iterate  $V_k$  is not necessarily the value function of some policy.

This algorithm is especially useful for discrete state and action spaces. It can be significantly sped up if we have sparse representations of the rewards and transitions.

The stopping condition for the iterations can either be that a fixed number of iterations  $K$  has been reached, or an  $\varepsilon$ -convergence criterion, i.e. stop when value function increments fall below some threshold (e.g.  $\|V_{k+1} - V_k\|_\infty \leq \varepsilon$ ).

**Proposition 2 (Convergence speed of VI)** *If given an upper bound  $r_{\max} \geq \sup |r(\cdot)|$  on the rewards, we  $\varepsilon$ -converge in at most*

$$K = \frac{\log(r_{\max})/\varepsilon}{\log(1/\gamma)}$$

*steps.*

**Performance guarantees.** The performance loss between the final policy  $\pi_K$  and the optimal policy  $\pi^*$  is measured by the distance between their values:

$$\|V^* - V^{\pi_K}\|_\infty \leq \frac{2\gamma}{1-\gamma} \|V^* - V_K\|_\infty$$

**Alternative termination condition.**

$$\|V_k - V_{k-1}\|_\infty - \min_{s \in \mathcal{S}} |V_k(s) - V_{k-1}(s)| \leq \varepsilon$$

**Complexity.** The temporal complexity of VI is  $O(KS^2A)$  – indeed the complexity of applying the Bellman operator is  $O(S^2A)$

## Variants

### 2.3.2 Policy Iteration

The idea of *policy iteration* is to use the action-state value function  $Q^\pi$  to iteratively update the policy.

**Greedy improvements.** We can improve over a policy  $\pi$  and define a new policy  $\pi'$  by acting greedily:

$$\pi'(s) := \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a)$$

It can then be shown that  $V^{\pi'} \geq V^\pi$ , and that improvements stop if and only if the Bellman optimality condition is satisfied.

To decide when to stop, we can use a set number of iterations, or use a  $\varepsilon$ -convergence criterion as suggested before for value iteration.

Policy iteration can also be performed for other improvement algorithms. For instance,  **$\varepsilon$ -greedy improvement** performs greedy updates at each state probabilistically, taking the greedy with probability  $1 - \varepsilon$  and otherwise taking a random action.



# Chapter 3

## Approximate Reinforcement Learning

As in regular control theory, dynamic programming opens up a lot of venues to solve MDPs we know everything about. But what do we do when we don't have a full, tabular model of the transitions and rewards of our system?

There are two ideas: use several **episodes** (trajectories) of agents following some policy (or policies) to compute empirical returns and build estimates of the value, or update our estimates as-we-go as we follow an agent's actions and the rewards it gets (**online learning**).

There are two aspects that need to be discussed:

- **policy evaluation**
- **policy optimization**

We will start by discussing approximate policy evaluation techniques and see how they apply to optimizing the policies as we go.

### 3.1 Monte Carlo Policy Evaluation

The idea is naïve: generate episodes  $\tau_i$  under a policy  $\pi$  starting at some state  $s_0$ , and compute their empirical returns

$$R(\tau_i) = \sum_{t=0}^{T_i} \gamma^t r_{t,i}.$$

The value estimate is then the empirical mean:

$$\hat{V}^\pi(s_0) = \frac{1}{n} \sum_{i=1}^n R(\tau_i) \quad (3.1)$$

Here, the beginning state  $s_0$  is fixed so we are only estimating its value.

**Dealing with non-episodic problems.** If the underlying problem does not necessarily terminate in finite time, we can always truncate it after a given number of time steps and say it has indeed “reset”. For instance, a trade execution algorithm does not have a terminal state so we might want to generate episodes by resetting after some time has passed. This means ignoring a term  $\sum_{t'=H+1}^{\infty} \gamma^{t'} r_{t'}$  in the return past a certain horizon  $H$ . In that case, the MC estimator converges to a truncated value function  $V_H^\pi$  which differs from the true value  $V^\pi$  by

$$|V_H^\pi(s_0) - V^\pi(s_0)| \leq \gamma^H \frac{\|r(\cdot)\|_\infty}{1 - \gamma}.$$

### 3.1.1 Incremental MC

We can easily see that the  $n$ -sample MC estimate can be seen as an update of the  $(n-1)$ -sample one:

$$\hat{V}_n^\pi(s_0) = \alpha_n R(\tau_n) + (1 - \alpha_n) \hat{V}_{n-1}^\pi(s_0)$$

where

$$\alpha_n = \frac{1}{n}$$

is the **learning rate** of the scheme.

Other learning rates can be used, and we have the following result ensuring convergence in the general case:

**Proposition 3** *Suppose the learning rate  $(\alpha_n)$  satisfies the Robbins-Monro condition:*

$$\sum_{n=0}^{\infty} \alpha_n = \infty \quad \sum_{n=0}^{\infty} \alpha_n^2 < \infty \quad (3.2)$$

*then the incremental MC estimate converges to the real value:*

$$\hat{V}^\pi(s_0) \xrightarrow{n \rightarrow \infty} V^\pi(s_0).$$

Incremental MC is also called TD(1). It does not allow online learning.

## 3.2 Temporal difference (TD) estimation

### 3.2.1 Using the Bellman equation: TD(0) estimation

The real value function  $V^\pi$  satisfies the Bellman equation. This means that the **temporal difference error** of a good estimate  $\hat{V}^\pi$  of  $V^\pi$ , defined after each transition and reward  $s_t, r_t, s_{t+1}$

$$\delta_t = r_t + \gamma \hat{V}^\pi(s_{t+1}) - \hat{V}^\pi(s_t), \quad (3.3)$$

should be small.

The TD(0) updates are performed as

$$\begin{aligned} \hat{V}^\pi(s_t) &\leftarrow \hat{V}^\pi(s_t) + \alpha_t \delta_t \\ &= (1 - \alpha_t) \hat{V}^\pi(s_t) + \alpha_t (r_t + \gamma \hat{V}^\pi(s_{t+1})) \end{aligned}$$

**Choosing the learning rate.** For incremental MC, a natural learning rate that works was  $\alpha_n = 1/n$ . For this temporal difference update scheme, it can be shown that we can use an **adaptive state-dependent learning rate**  $\alpha_t = \alpha(s_t)$ :

**Proposition 4** *Set the learning rate to  $\alpha_t = \alpha(N_t(s))$  where  $N_t(s)$  is the number of times state  $s \in \mathcal{S}$  has been visited before time  $t$ . If we suppose  $N_t(s) \rightarrow \infty$  (every state is visited infinitely many times), then the conditions*

$$\sum_{n=0}^{\infty} \alpha(n) = \infty \quad \sum_{n=0}^{\infty} \alpha(n)^2 < \infty \quad (3.4)$$

lead to convergence of the scheme.

This scheme doesn't require the entire trajectory to be known to perform an update, allowing for fully online learning.

### 3.2.2 Temporal difference TD( $\lambda$ )

These schemes average between incremental MC and TD(1). The scheme is defined for  $0 < \lambda < 1$  as:

$$\hat{V}^\pi(s_t) \leftarrow \hat{V}^\pi(s_t) + \alpha(s_t) \sum_{t'=t}^T (\gamma\lambda)^{t'-t} \delta_{t'} \quad (3.5)$$

where  $(\delta_{t'})_{t'}$  is the sequence of temporal differences errors of the trajectory.

This scheme *does* require the entire trajectory to be known before performing updates, but it re-discounts them using the parameter  $\lambda$ .

#### Eligibility traces

At every transition  $s_t, r_t, s_{t+1}$ , we define the temporal difference  $\delta_t$  and update the **eligibility traces**  $z(s)$  as

$$\forall s \in \mathcal{S} \quad z(s) = \begin{cases} \lambda z(s) & \text{if } s \neq s_t \\ 1 + \lambda z(s) & \text{if } s = s_t \\ 0 & \text{if } s_t = s_0 \text{ (reset the traces)} \end{cases}$$

For policy evaluation, we update the value function estimate for all states  $s \in \mathcal{S}$  using the eligibility traces with

$$\hat{V}^\pi(s) \leftarrow \hat{V}^\pi(s) + \alpha(s) z(s) \delta_t \quad (3.6)$$

This introduces a compromise between propagating rewards faster than TD(0) and having smaller variance than incremental MC.

**Remark 4** All the previous ideas for policy evaluation on the value function  $V^\pi$  readily apply to evaluating the action-state value  $Q^\pi$ .

## 3.3 Application of TD to policy improvement

After introducing the previous temporal difference notions for policy evaluation, we can see how they work for **policy improvement** schemes. The basic idea is to alternate the previous temporal-difference policy evaluation techniques with policy improvement ideas from chapter 2 which exploit the  $Q$  value function.

### 3.3.1 SARSA

We define an exploration policy using the softmax function and a temperature parameter  $\tau$ :

$$\pi_Q(s, a) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a' \in \mathcal{A}} \exp(Q(s, a')/\tau)}$$

After taking an action according to  $\pi_Q$ , we update using the TD(1) scheme: compute the temporal difference of the  $Q$ -function

$$\delta_t = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t) \quad (3.7)$$

and update  $Q$

$$\begin{aligned}\widehat{Q}(s_t, a_t) &\leftarrow \widehat{Q}(s_t, a_t) + \alpha_t \delta_t \\ &= (1 - \alpha_t) \widehat{Q}(s_t, a_t) + \alpha_t (r_t + \widehat{Q}(s_{t+1}, a_{t+1}))\end{aligned}$$

We can ensure the policy continues to improve by decreasing the temperature  $\tau$  to 0 and becoming more and more greedy in our action selection (keeping the same  $\tau$  reduces the algorithm to some kind of softmax policy evaluation).

Extension to TD( $\lambda$ ) update schemes is straightforward.

### 3.3.2 Q-learning

The idea is to introduce the optimal TD error for the  $Q$  value function

$$\delta_t = r_t + \gamma \max_{a'} \widehat{Q}(s_{t+1}, a') - \widehat{Q}(s_t, a_t) \quad (3.8)$$

and updating the  $Q$  estimate using that:

$$\widehat{Q}(s_t, a_t) \leftarrow \widehat{Q}(s_t, a_t) + \alpha_t \delta_t$$

**Remark 5** This is *not* the same as taking the greedy action when simulating the trajectory: this is upper-bounding the  $Q$ -value TD error (3.7) and using that upper bound to update the value.

---

#### Algorithm 3.1: Q-learning

---

**Input:** Initial state  $s_0$ , exploration policy

```

1 repeat
2   Take action  $a_t$  according the exploration policy;
3   Get reward and next state  $(r_t, s_{t+1})$ ;
4   Compute TD error  $\delta_t = r_t + \gamma \max_{a'} \widehat{Q}(s_{t+1}, a') - \widehat{Q}(s_t, a_t)$ ;
5    $\widehat{Q}(s_t, a_t) \leftarrow \widehat{Q}(s_t, a_t) + \alpha_t \delta_t$ ;
6 until trajectory terminates;
   // Compute the policy from the  $Q$  value
7 foreach  $s \in \mathcal{S}$  do
8    $\tilde{\pi}(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \widehat{Q}(s, a)$ ;
9 return Policy  $\tilde{\pi}$ ;
```

---

We can use several exploration policies, such as  $\varepsilon$ -greedy exploration (where we take the  $\widehat{Q}$ -optimal action with probability  $\varepsilon$  and otherwise take a random action).

**What about the learning rate.** We can introduce, just as in evaluating the value  $V$  (see proposition 4), a action-state count  $N_t(s, a)$  and set an adaptive learning rate

$$\alpha_t = \alpha(N_t(s_t, a_t))$$

and the Robbins-Monro conditions on the function  $\alpha$  work out the same.

## Chapter 4

# Approximate solving of Markov Decision Processes with Policy Gradients

Solving MDPs is seeking the maximizing policy of the value function. For approximate solving of MDPs, we target what could be a more general **policy performance metric**. Often, it is indeed the value function

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^T \gamma^t r_t \right] = \mathbb{E}_{\tau \sim \pi} [R(\tau)] \quad (4.1)$$

where  $\tau = \{s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T\}$  and  $R(\tau) = \sum_{t=0}^T \gamma^t r_t$  is the total return of the trajectory  $\tau$ . We use the shorthand  $r_t$  for the reward  $r(s_t, a_t)$  – or more generally a transition reward  $r(s_t, a_t, s_{t+1})$  or stochastic reward distributed as  $r_t \sim p(\cdot, s_t, a_t, s_{t+1})$ . The expectation  $J$  is either conditional on a given starting point  $s_0$ , or on a distribution for it.<sup>1</sup>

We seek to compute the maximizing policy in a parametric search space  $\{\pi_\theta : \theta \in \Theta\}$ :

$$\max_{\theta} J(\pi_\theta)$$

Iterative methods could be used if  $J$  can be obtained in closed form; if these are not given they could be estimated by Monte Carlo methods, which would be very expensive and wasteful.

Instead, we will update the estimate as we go and simulate trajectories, by iteratively updating the policy parameter  $\theta$  using a gradient ascent method with an estimated gradient.

**Proposition 5 (Gradient under a parametric law)** *Given a set of probability models  $\{P_\theta : \theta \in \Theta \subseteq \mathbb{R}^d\}$  on a set  $\mathcal{X}$  and a function  $f: \mathcal{X} \rightarrow \mathbb{R}$ , we have that*

$$\nabla_{\theta} \mathbb{E}_{X \sim P_\theta} [f(X)] = \mathbb{E}_{X \sim P_\theta} [f(X) \nabla_{\theta} \log P_\theta(X)]$$

*This is a useful property for deriving estimators of the derivatives in optimization problems with stochastic objectives.*

*Generalization to the case where  $f$  also depends on  $\theta$  is straightforward.*

This can be shown either by either writing the expectation as an integral, or by a change of measures with a Radon-Nikodym derivative.

Proposition 5 allows us to write the gradient of (4.1), called the **policy gradient** as an expectation:

$$\nabla_{\theta} J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ R(\tau) \sum_{t=0}^T \nabla_{\theta} \log \pi_\theta(s_t, a_t) \right] \quad (4.2)$$

<sup>1</sup>For instance, OpenAI Gym's **CartPole-v1** environment has a stochastic initial state  $s_0$ .

and we will need to derive estimations for this quantity.

There are other ways of writing the policy gradient, such as (see [1, chap. 13])

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_s \left[ \sum_a Q^{\pi_{\theta}}(s, a) \nabla \pi_{\theta}(s, a) \right]$$

## 4.1 Monte Carlo policy gradient: the REINFORCE algorithm

**The idea.** The policy gradient (4.2) is an expectation which can be estimated using Monte Carlo approximation. We obtain the following estimate:

$$\widehat{\nabla_{\theta} J(\pi_{\theta})} = \frac{1}{M} \sum_{i=1}^M R(\tau_i) \sum_{t=0}^{T_i} \nabla_{\theta} \log \pi_{\theta}(s_t^i, a_t^i) \quad (4.3)$$

This is an unbiased Monte Carlo estimate of the policy gradient. It only requires suitable regularity of the parametric policy model  $\theta \mapsto \pi_{\theta}$ .

**Remark 6** Equation (4.3) can be used as-is for functions with simple closed-form derivatives. In an automatic differentiation framework such as PyTorch, we can instead get the policy gradient from a computational graph with the following pseudo-loss function:

$$\begin{aligned} \tilde{J}(\theta) &= \frac{1}{M} \sum_{i=1}^M R(\tau_i) \sum_{t=0}^{T_i} \log \pi_{\theta}(s_t^i, a_t^i) \\ &= \frac{1}{M} \sum_{i=1}^M \left( \sum_{t=0}^{T_i} \gamma^t r_t^i \right) \sum_{t=0}^{T_i} \log \pi_{\theta}(s_t^i, a_t^i) \end{aligned} \quad (4.4)$$

We will simulate multiple trajectories (*episodes*) to perform the gradient update: this is the similar to what is done with batch, mini-batch or stochastic gradient steps. This leads to algorithm 4.1.

---

### Algorithm 4.1: Monte Carlo Policy Gradient (REINFORCE)

---

**Input:** Arbitrary initial policy  $\pi_{\theta_0}$ .

**Output:** Optimal parametric policy  $\pi_{\theta^*}$ .

```

1 repeat
2   Simulate a trajectory  $\tau$ ;
3    $\mathbf{g} \leftarrow (\sum_{t=0}^T \gamma^t r_t) \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(s_t, a_t);$            // policy gradient
4    $\theta \leftarrow \theta + \alpha \mathbf{g};$                                            // update parameter
5 until finished;
```

---

This algorithm can be modified in several ways, by performing the gradient step at every time in the process, or on the entire batch of trajectories. Sutton and Barto [1] actually use the more fine-grained update scheme wherein the parameter  $\theta$  is updated at each step  $t$  of every episode<sup>2</sup>.

### 4.1.1 Variance reduction: temporal structure

We can re-weight the log-probability gradients in eq. (4.2) by exploiting the fact that, for any time  $t$ , the cumulative rewards  $\sum_{t'=0}^{t-1} \gamma^{t'} r_{t'}$  from 0 to  $t-1$  are measurable with respect to the trajectory up to  $t$ ,  $\tau_{0:t}$ :

---

<sup>2</sup>This allows for fully online learning.

**Proposition 6** The policy gradient (4.2) can be rewritten as

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \sum_{t'=t}^T \gamma^{t'} r_{t'} \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \right] \quad (4.5)$$

which leads to the policy gradient estimate

$$\widehat{\nabla_{\theta} J}(\pi_{\theta}) = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^{T_i} \gamma^t G_t^i \nabla_{\theta} \log \pi_{\theta}(s_t^i, a_t^i) \quad (4.6)$$

where<sup>a</sup>  $G_t^i = \sum_{t'=t}^{T_i} \gamma^{t'-t} r_{t'}^i$ .

<sup>a</sup>This quantity is an estimate of the  $Q$ -function  $Q^{\pi}(s_t, a_t) = \mathbb{E}[\sum_{t'=t}^T \gamma^{t'-t} r_{t'} \mid s_t, a_t]$ .

The following algorithm provides an efficient recursive method for computing the returns  $G_t^i$ :

---

**Algorithm 4.2:** Computing the returns

---

**Input:** The rewards  $(r_t)_{0 \leq t \leq T}$ , discount factor  $\gamma$

**Output:** The array of discounted returns.

```

1  $R_T \leftarrow r_T$ ;
2 foreach  $t \leftarrow T - 1$  to 0 do
3    $G_t \leftarrow r_t + \gamma G_{t+1}$ ;
4 return  $(G_t)_{0 \leq t \leq T}$ ;
```

---

### 4.1.2 Variance reduction with baselines

Given any **baseline** function  $b: \mathcal{S} \rightarrow \mathbb{R}$ , we can rewrite the policy gradient again as

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \left( \sum_{t'=t}^T \gamma^{t'} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \right] \quad (4.7)$$

The resulting policy gradient estimate we get is

$$\widehat{\nabla_{\theta} J}(\pi_{\theta}) = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^{T_i} \gamma^t (G_t^i - b(s_t^i)) \nabla_{\theta} \log \pi_{\theta}(s_t^i, a_t^i) \quad (4.8)$$

which is still an unbiased estimate – but with the added benefit of **variance reduction** if  $b$  is chosen right.

#### The best baseline.

It can be shown that the “best” baseline in terms of variance reduction is actually the value function:

$$b^*(s_t) = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \mid s_t \right] = V^{\pi}(s_t)$$

...which we are already trying to approximate. This suggests that we do some kind of **bootstrapping**. The approximate policy evaluation ideas from the previous chapter 3 will come in handy.

The obvious first idea is nonparametric Monte Carlo estimation:

$$\widehat{b}(s_t) = \frac{1}{M} \sum_{i=1}^M G_t^i$$

This baseline is an unbiased estimate: following eq. (2.6) the expectation of the  $Q$ -function estimate  $G_t$  under the policy  $\pi_\theta$  is

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi_\theta(s_t, \cdot)} \mathbb{E}_{\pi_\theta} [G_t] = \mathbb{E}_{a_t \sim \pi_\theta(s_t, \cdot)} \mathbb{E}_{\pi_\theta} \left[ \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \mid s_t, a_t \right]$$

However, this introduces additional variance and the baseline might become noisy.

### 4.1.3 Parametric Bootstrapping of the baseline

We define the bootstrap estimate  $b(s) := \widehat{v}_\nu(s)$  lying in a parametric search space  $\{\widehat{v}_\nu : \nu \in \mathcal{V}\}$ . The value parameter can be iteratively updated using gradient steps alternatively with the policy parameter  $\theta$ .

For a given trajectory sample  $\tau = \{s_0, a_0, r_0, \dots\}$ , introduce the mean-squared error between the forward returns (a nonparametric estimate of the value function) and the value model's output:

$$\mathcal{L}(\nu; \tau) = \sum_{t=0}^T (G_t - \widehat{v}_\nu(s_t))^2 \quad (4.9)$$

Then before each update of the policy  $\pi_\theta$ , update the value parameter  $\nu$  using either the gradient of  $\mathcal{L}$ .

The adapted episodic learning algorithm with an adaptive baseline is as follows:

---

**Algorithm 4.3:** REINFORCE with parametric baseline

---

```

1 repeat
2   Simulate a trajectory  $\tau$ ;
3   Compute the returns  $(G_t)$  of the trajectory;
4    $\mathbf{g} \leftarrow \sum_{t=0}^T \gamma^t (G_t - \widehat{v}_\nu(s_t)) \nabla_\theta \log \pi_\theta(s_t, a_t)$  ;           // policy gradient
5    $\nu \leftarrow \nu - \beta \nabla_\nu \mathcal{L}(\nu; \tau)$  ;                                   // update value estimate
6    $\theta \leftarrow \theta + \alpha \mathbf{g}$  ;                                           // update policy
7 until finished;
```

---

As before (see remark 6), this algorithm can be implemented within an automatic differentiation framework such as PyTorch by defining the right computational graphs. The associated pseudo-loss would be

$$\tilde{J}(\theta) = \sum_{t=0}^T \gamma^t (G_t - \widehat{v}_\nu(s_t)) \log \pi_\theta(s_t, a_t)$$

**Remark 7** *The difference between the returns and the value approximation  $G_t - \widehat{v}_\nu(s_t)$  is a (biased) estimate of the temporal-difference error  $\delta_t = r_t + \gamma \widehat{v}(s_{t+1}) - \widehat{v}(s_t)$ . Thus, the mean-squared error  $\mathcal{L}$  can be seen as an estimate of the Bellman error of the value proposal  $\widehat{v}_\nu$ .*



## 4.2 Bootstrapped returns: Actor-Critic algorithms

**The idea.** The enhanced REINFORCE algorithm builds estimates of the value function to compute the policy gradient: this Monte Carlo method is computationally expensive, and may still lead to high variance. To combat this, it might be a good idea to *learn* from the Monte Carlo estimates in a way that gives a consistent estimate that follows the policy gradient updates.

To achieve this, the class of **actor-critic methods** introduces a second search space for approximation of the state(-action) value function and uses **bootstrapped temporal-difference estimates** for the returns  $G_t$ , by taking inspiration from TD policy evaluation (see chapter 3). The supervised baseline of section 4.1.3 was already a first step towards this, but this time we actually use the value estimate for policy evaluation and not only for variance reduction. The subtle difference is further explained in [1, chap. 13.5].

### 4.2.1 Actor-Critic

Learning the policy is still done by gradient steps, using estimates of the form eq. (4.8). But this time, we replace in the objective  $J$  the Monte Carlo returns  $G_t$ , which estimate the  $Q$ -function, by a parametric estimator  $\hat{q}_\omega(s_t, a_t)$  called the **critic**. As the algorithm runs, the critic learns how to value each action suggested by the learned policy  $\pi$  (called the **actor**) by looking at the bootstrapped return estimates. The policy gradient becomes

$$\widehat{\nabla_\theta J(\pi_\theta)} = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^{T_i} \hat{q}_\omega(s_t^i, a_t^i) \nabla_\theta \log \pi_\theta(s_t^i, a_t^i) \quad (4.10)$$

Now, we want the critic  $\hat{q}_\omega$  to minimize the **Bellman error**

$$\mathbb{E}_{\pi_\theta} [(r_t + \gamma \hat{q}_\omega(s_{t+1}, a_{t+1}) - \hat{q}_\omega(s_t, a_t))^2 \mid s_t, a_t]$$

because having it as small as possible will ensure that it follows the dynamic programming principle, making it a “good” estimate of the real action-state value function  $Q^{\pi_\theta}$ . The above expectation can be estimated using the TD(0) error of  $\hat{q}_\omega$ :

$$\delta_t = r_t + \gamma \hat{q}_\omega(s_{t+1}, a_{t+1}) - \hat{q}_\omega(s_t, a_t).$$

which we will seek to minimize for every time step or episode alternatively with the actor loss.

This leads to the following algorithm:

---

#### Algorithm 4.4: Actor-Critic

---

**Input:** Initial policy parameter  $\theta$ , value parameter  $\omega$

**Output:** Policy  $\pi_{\theta^*}$ , action-state value approximation  $\hat{q}_{\omega^*}$

---

```

1 repeat
2   Simulate trajectory  $\tau$ ;
3   foreach  $t = 0, \dots, T$  do
4      $\delta_t \leftarrow r_t + \gamma \hat{q}_\omega(s_{t+1}, a_{t+1}) - \hat{q}_\omega(s_t, a_t)$  ;           // TD(0) error
5      $\omega \leftarrow \omega + \beta \sum_{t=0}^T \delta_t \nabla_\omega \hat{q}_\omega(s_t, a_t)$  ;           // critic update
6      $\theta \leftarrow \theta + \alpha \sum_{t=0}^T \hat{q}_\omega(s_t, a_t) \nabla_\theta \log \pi_\theta(s_t, a_t)$  ; // actor update
7 until finished;
```

---

### 4.2.2 Actor-Critic with baselines: Advantage Actor-Critic (A2C)

As we've seen before, a good baseline to introduce in the policy gradient is the value function – see eq. (4.7). The coefficient before the log-probability gradient in (4.2) becomes an estimate of  $Q^\pi(s, a) - V^\pi(s) = \mathbb{E}_{s'}[r(s, a) + \gamma V^\pi(s') \mid s, a]$ . We introduce the **advantage** function

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s). \quad (4.11)$$

An unbiased estimate of  $A^\pi$  is the temporal difference error of the value function

$$\delta_t = r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t).$$

This can be seen by taking the expectation with respect to  $s_{t+1}$  conditionally on  $(s_t, a_t)$ .

This leads to the (bootstrapped) estimate, using a parametric approximation  $\hat{v}_\nu$  of the state value function:

$$\hat{\delta}_t = r_t + \gamma \hat{v}_\nu(s_{t+1}) - \hat{v}_\nu(s_t).$$

Now, we have no need to learn a action-state critic  $\hat{q}_\omega$ , since we have TD error estimates that can be expressed only using  $\hat{v}_\nu(s)$  – which is simpler.

The A2C algorithm is as follows:

---

**Algorithm 4.5:** Advantage Actor-Critic (A2C)

---

**Input:** Initial policy parameter  $\theta$ , value parameter  $\nu$

**Output:** Policy  $\pi_{\theta^*}$ , value approximation  $\hat{v}_{\nu^*}$

---

```

1 repeat
2   Simulate trajectory  $\tau$ ;
3   foreach  $t = 0, \dots, T$  do
4      $\delta_t \leftarrow r_t + \gamma \hat{v}_\nu(s_{t+1}) - \hat{v}_\nu(s_t)$  ; // TD(0) error
5      $\nu \leftarrow \nu + \beta \sum_{t=0}^T \delta_t \nabla_\nu \hat{v}_\nu(s_t)$  ; // critic update
6      $\theta \leftarrow \theta + \alpha \sum_{t=0}^T \hat{v}_\nu(s_t) \nabla_\theta \log \pi_\theta(s_t, a_t)$  ; // actor update
7 until finished;
```

---

**Remark 8** To leverage automatic differentiation, appropriate pseudo-losses to define a computational graph for the updates in algorithm 4.5 are

$$\tilde{C}(\nu) = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^T \delta_t^i \hat{v}_\nu(s_t^i) \quad (4.12a)$$

$$\tilde{J}(\theta) = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^T \hat{v}_\nu(s_t^i) \log \pi_\theta(s_t^i, a_t^i) \quad (4.12b)$$

where the temporal-difference estimates  $\delta_t^i$  must be detached from the graph.

The value update step in algorithm 4.5 (and also algorithm 4.4) can be seen as a supervised regression step, where we fit the (bootstrapped) returns  $G_{t:t+1} = r_t + \gamma \hat{v}_\nu(s_{t+1})$  to the value model predictions  $\hat{v}_\nu(s_t)$ , with a mean-squared loss

$$\mathcal{L}(\nu) = \sum_{t=0}^T (G_{t:t+1} - \hat{v}_\nu(s_t))^2 \quad (4.13)$$

We then perform a semi-gradient update step, taking the gradient with respect to  $\nu$  whilst ignoring the dependency of the return estimate  $G_{t:t+1}$  on it.

**Remark 9** This means that the pseudo-loss  $\tilde{C}(\nu)$  in remark 8 can be replaced by the MSE  $\mathcal{L}(\nu)$ .

We create the estimate  $G_{t:t+1}$  by only looking at the immediate reward and bootstrapping the rest of the value: this is called **one-step actor-critic**. It has poor sample efficiency: the return estimates ignore most of the (actual) future rewards and end up bootstrapping too much. However, it naturally allows for *online learning*.

This can be easily generalized to using more rewards for bootstrapping the returns:  **$n$ -step actor-critic** methods look further ahead and use estimates

$$G_{t:t+n} = \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n \hat{v}_\nu(s_{t+n})$$

**Batch actor-critic.** In batch algorithms, we switch out episodic simulation with epochs of simulation-and-reset batches, where we simulate trajectories (and reset them at terminal states) until a fixed number  $B$  of updates have been done.

---

**Algorithm 4.6:** Batch A2C

---

**Input:** Number of epochs  $N$ , batch size  $B$

```

1 foreach  $n = 1, \dots, N$  do
2    $s_0 \leftarrow \text{Reset}()$ ;
3   foreach  $i = 0, \dots, B - 1$  do
4      $v_i \leftarrow \hat{v}_\nu(s_i)$ ; // estimate value
5     Draw action  $a_i \sim \pi_\theta$ ;
6     Get state, reward  $(r_i, s_{i+1})$ ;
7     // check if we just finished a trajectory
8     if  $s_{i+1}$  is terminal then
9        $s_{i+1} \leftarrow \text{Reset}()$ ;
10    if  $s_B$  is terminal then
11       $v_B \leftarrow 0$ ;
12    else
13       $v_B \leftarrow \hat{v}_\nu(s_B)$ ;
14    Compute bootstrapped advantages, returns  $(\delta_i, \hat{G}_i)$ ;
15     $\nu \leftarrow \nu + \beta \sum_{i=0}^{B-1} \delta_i \nabla_\nu \hat{v}_\nu(s_i)$ ;
16     $\theta \leftarrow \theta + \alpha \sum_{i=0}^{B-1} \delta_i \nabla_\theta \log \pi_\theta(s_i, a_i)$ ;
```

---

The returns  $\hat{G}_i$  have to be computed taking into account when trajectories end. An efficient algorithm is

$$\hat{G}_i = \begin{cases} r_i + \gamma \hat{G}_{i+1} & \text{if } s_i \text{ is not terminal} \\ r_i & \text{otherwise} \end{cases}$$

for  $i = 0, \dots, B - 2$  and  $G_{B-1} = r_{B-1} + \gamma v_B$ . The TD(0) errors are  $\delta_i = \hat{G}_i - v_i$ .

# **Chapter 5**

## **The exploration-exploitation dilemma**

### **5.1 Multi-Armed Bandits**

# Bibliography

- [1] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. Second. The MIT Press, 2018. URL: <http://incompleteideas.net/book/the-book-2nd.html>.