

# CRASH COURSE IN REINFORCEMENT LEARNING

Damian Leśniak  
Aleksandra Nowak  
Igor Sieradzki  
Maciej Wołczyk

November 25, 2019



# PART 1: INTRODUCTION TO REINFORCEMENT LEARNING

# REINFORCEMENT LEARNING



Images (from upper left): <https://waymo.com/press/>, <https://starcraft2.com/pl-pl/media>,  
<https://cloud.google.com/tpu/>, <https://ai.googleblog.com/2018/06/scalable-deep-reinforcement-learning.html>

# AGENT AND ENVIRONMENT

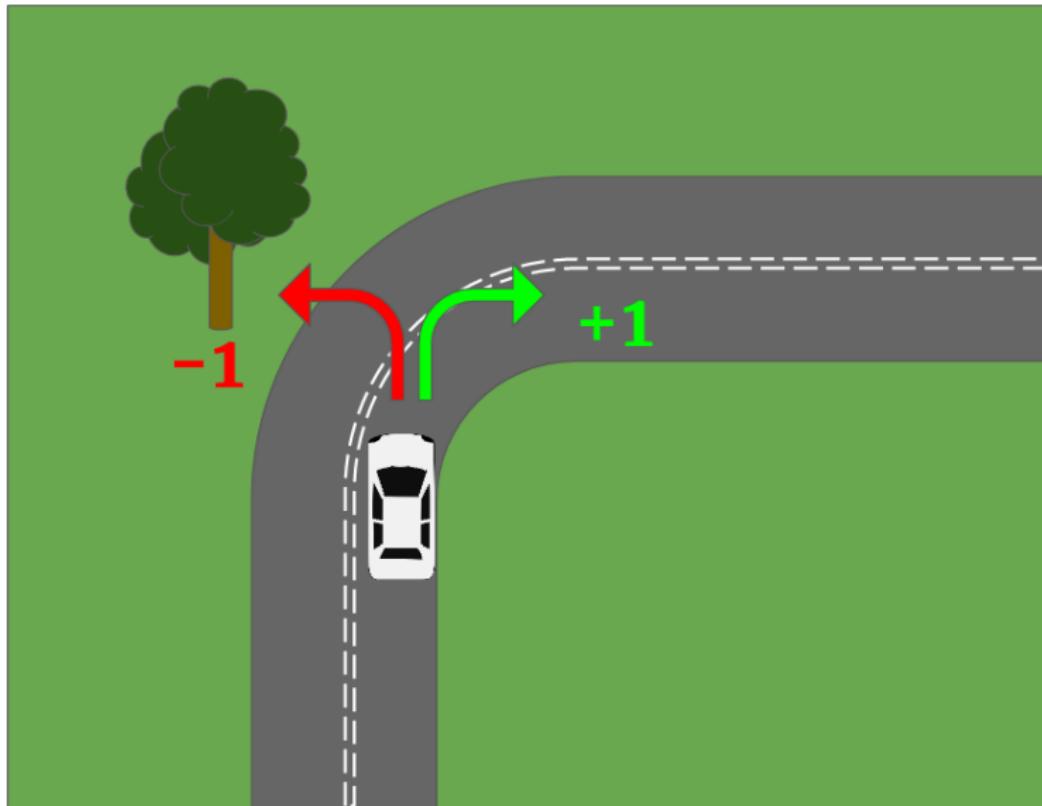


# ACTIONS AND OBSERVATIONS

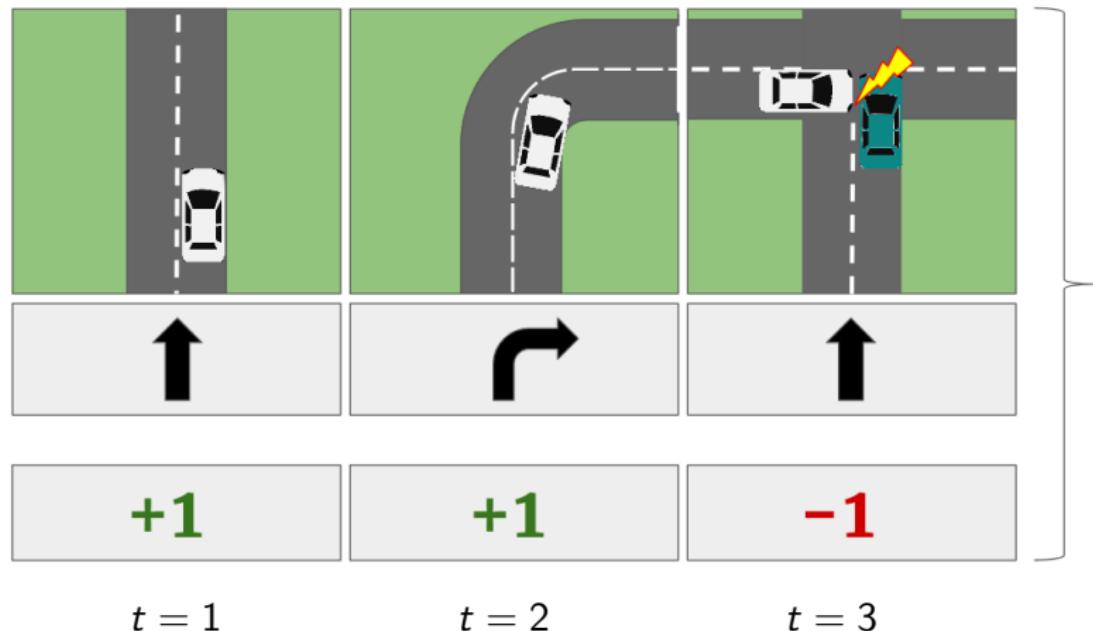


Images: <https://waymo.com/press/>, <https://research.mapillary.com/>

## REWARDS



# TRAJECTORIES



## NOTATION

1.  $o \in \mathbb{O}$  is the observation (vector of real numbers), describing what the agent sees.
2.  $a \in \mathbb{A}$  is the taken action.
3.  $r \in \mathbb{R}$  is the reward for performing action  $a$ .
4. For simplicity, let us assume that the action space is finite and has  $A$  elements

$$\mathbb{A} = \{a^{(1)}, \dots, a^{(A)}\}$$

## NOTATION

To distinguish between observations, actions and rewards obtained at different (discrete) time  $t$ , we write:  $o_t \in \mathbb{O}$ ,  $a_t \in \mathbb{A}$ ,  $r_t \in \mathbb{R}$ .

Note the difference between the lower (time) and upper (indicator of the action) indices. In particular:

- $a_3 = a_4$  – action taken at time  $t = 3$  is equal to the one taken at time  $t = 4$
- $a_4 = a^{(7)}$  – action taken at time  $t = 4$  is the seventh possible action (for example:  $a^{(7)} = \text{turn left}$ ).
- $a_7 = a^{(1)}$  – action taken at time  $t = 7$  is equal to the first possible action (for example  $a^{(1)} = \text{go forward}$ ). Note that  $a_7 \neq a^{(7)}$ .

# TRAJECTORIES

Possible actions:

$$\mathbb{A} = \{a^{(1)} = \text{Left } \leftarrow, a^{(2)} = \text{Forward } \uparrow, a^{(3)} = \text{Right } \rightarrow\}$$



+1

+1

-1

$$a_1 = a^{(2)}$$

$$a_2 = a^{(3)}$$

$$a_3 = a^{(2)}$$

## NOTATION – AGENT

At time  $t$ , the agent is given the observation vector  $o_t$ , based on which it performs an action  $a_t$  and receives the reward  $r_t$ .

The agent's behaviour for each  $o \in \mathbb{O}$  is its *policy*. The agent's policy defines a conditional probability over the possible choices for action  $a$  given observation  $o$ .

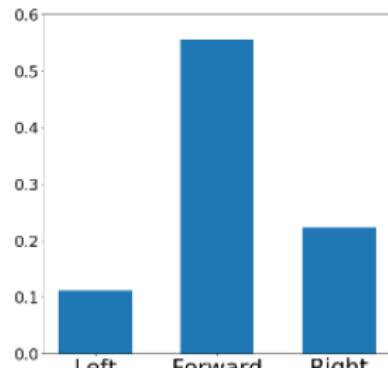
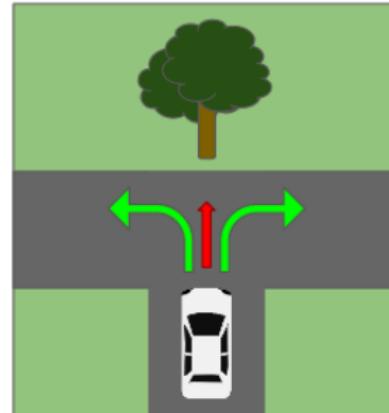
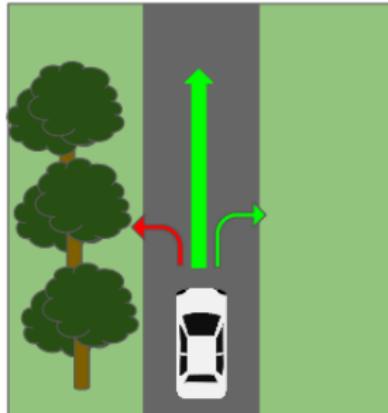
Probability distributions on finite sets can be conveniently represented as vectors, therefore we write:

$$\pi : \mathbb{O} \ni o \mapsto (\pi(a^{(1)} | o), \dots, \pi(a^{(A)} | o)) \in \mathbb{R}^A$$

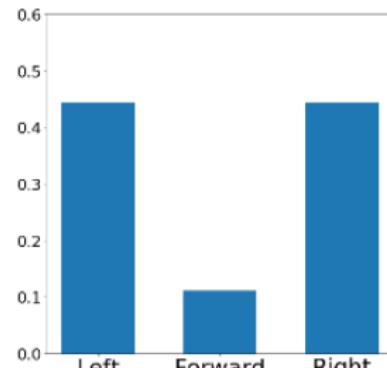
$$0 \leq \pi(a^{(1)} | o), \dots, \pi(a^{(A)} | o) \leq 1$$

$$\pi(a^{(1)} | o) + \dots + \pi(a^{(A)} | o) = 1$$

## NOTATION – AGENT



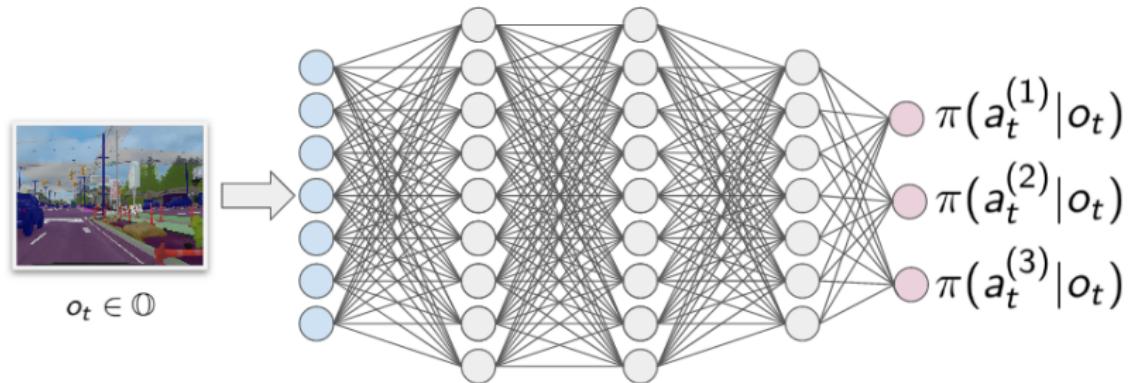
$$\pi(a_1|o_1)$$



$$\pi(a_2|o_2)$$

## NOTATION – AGENT: A NEURAL NETWORK

The agent may be modeled by a neural network, which defines its policy  $\pi_\theta$  with respect to trainable parameters  $\theta$ :



Note: we need to assign each action its positional number from the set  $\{1, 2, \dots, A\}$ , in order to identify the outputs from the network with actions – similar to classification tasks and labels!

## NOTATION – ENVIRONMENT

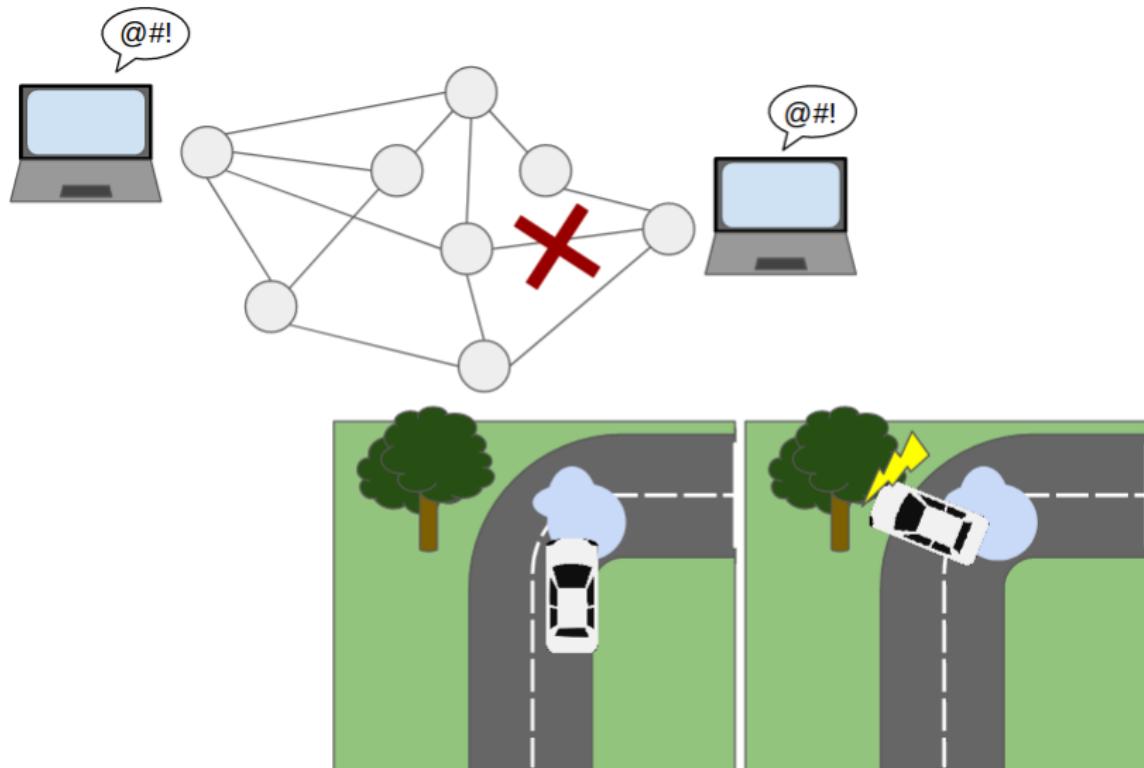
### ENVIRONMENT

The environment is the overall model of the world. At any given time step  $t$  it is described by a state  $s_t \in \mathbb{S}$ . After the agent's action, the environment:

1. Emits the rewards for the taken action  $r_t$ .
2. Changes the state to  $s_{t+1}$ .
3. Emits observations  $o_{t+1}$ .

From now on we assume that  $o_t = s_t$  for all  $t$ . This means that in each step the agent sees the full state of the environment.  
Hereinafter we won't be using the  $o_t$  notation anymore.

# ENVIRONMENTS MAY NOT BE DETERMINISTIC



## ENVIRONMENT – MARKOV PROPERTY

### MARKOV PROPERTY

The environment changes its state to  $s_{t+1}$  accordingly to the *conditional probability distribution* given the current state and the agent action. Note that the future state does not depend on the past, but only on the present – i.e. the state transition satisfies the Markov property:

$$(r_t, s_{t+1}) \sim p(r_t, s_{t+1} | s_t, a_t, r_{t-1}, s_{t-1}, a_{t-1}, \dots, s_1) = p(r_t, s_{t+1} | s_t, a_t)$$

## NOTATION – PROBABILITY DISTRIBUTIONS

All probabilities that (even in part) depend on  $\theta$  will be denoted by  $\pi_\theta$ , other by  $p$ .

- $p(s_1)$  – initial state (environment)
- $\pi_\theta(a_t | s_t)$  – action (agent)
- $p(r_t, s_{t+1} | s_t, a_t)$  – state transition (environment)
- $p(r_t, s_{\text{END}} | s_t, a_t)$  – end of episode (environment)
- $\pi_\theta(\tau)$  – trajectory (agent + environment)

## NOTATION – TRAJECTORY

A trajectory is a series of pairs of states and taken actions. In this presentation we also include rewards. It's more intuitive and corresponds to the implementation we will use.

### TRAJECTORY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

### THE PROBABILITY OF TRAJECTORY $\tau$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

## TRAJECTORY – PROBABILITY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

## TRAJECTORY – PROBABILITY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_1, a_1, r_1, s_2, \dots, s_T, a_T)$$

## TRAJECTORY – PROBABILITY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_1, a_1, r_1, s_2, \dots, s_T, a_T)$$

$$\stackrel{\text{Markov}}{=} \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_T, a_T)$$

## TRAJECTORY – PROBABILITY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\begin{aligned}\pi_\theta(\tau) &= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}}) \\ &= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_1, a_1, r_1, s_2, \dots, s_T, a_T) \\ &\stackrel{\text{Markov}}{=} \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_T, a_T) \\ &= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T) \ \pi_\theta(a_T | s_{T-1}) \ p(r_T, s_{\text{END}} | s_T, a_T)\end{aligned}$$

# TRAJECTORY – PROBABILITY

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_1, a_1, r_1, s_2, \dots, s_T, a_T)$$

$$\stackrel{\text{Markov}}{=} \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T, a_T) \ p(r_T, s_{\text{END}} | s_T, a_T)$$

$$= \pi_\theta(s_1, a_1, r_1, s_2, \dots, s_T) \ \pi_\theta(a_T | s_{T-1}) \ p(r_T, s_{\text{END}} | s_T, a_T)$$

= ...

$$= p(s_1) \ \pi_\theta(a_1 | s_1) \ p(r_1, s_2 | s_1, a_1) \dots p(r_T, s_{\text{END}} | s_T, a_T)$$

## ENVIRONMENT – REWARDS

For every action  $a_t$  the environment returns a reward  $r_t$ .

### PROBLEM: WHAT IS THE REWARD?

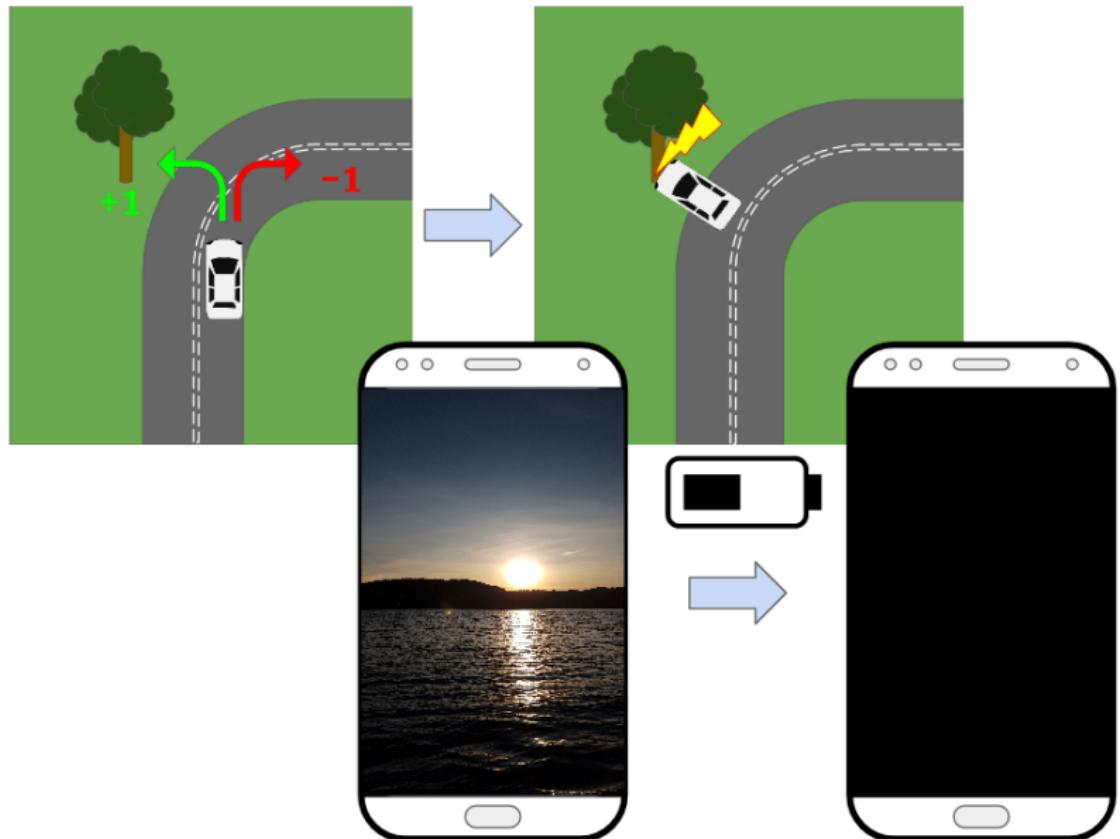
The real (i.e. ideal) reward distribution is usually unknown and task-dependent. In practice, the rewards are defined by the programmer.

While designing the rewards, one has to take into account all possible outcomes of each situation and assign them a real number. It is usually very unclear what that number should be.

## ENVIRONMENT – BAD REWARDS

One of the key challenges in developing a Reinforcement Learning agent is to correctly design the rewards. Note that bad rewards may lead to wrong or unexpected behaviours.

## ENVIRONMENT – BAD REWARDS



## RETURN

The return  $R_t$  at time  $t$  is the sum of all the rewards starting from time  $t$  to the end of episode  $T$ , obtained by following a given policy  $\pi_\theta$ :

$$R_t = r_t + r_{t+1} + r_{t+2} + \cdots + r_T = \sum_{j=t}^T r_j$$

Therefore the return of trajectory  $\tau$  is simply:

$$R(\tau) = R(s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}}) = \sum_{j=1}^T r_j$$

## PART 2: POLICY GRADIENT

## LEARNING OBJECTIVE

The learning objective is to **maximize** the expected return

$$\mathbb{E}_{\tau \sim \pi_\theta} R(\tau) \approx \frac{1}{N} \sum_{j=1}^N R(\tau_j),$$

where  $\{\tau_1, \dots, \tau_N\}$  is the training set of  $N$  trajectories generated by policy  $\pi_\theta$ .

## GRADIENT

The maximization of the above function (of learnable parameters  $\theta$ ) with gradient-based methods requires the computation of the gradient:

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} R(\tau)$$

Then we may update the parameters  $\theta$  by taking an ascending step in the direction of the gradient.

## GRADIENT

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} R(\tau) &= \nabla_{\theta} \int \pi_{\theta}(\tau) R(\tau) d\tau \\&= \int \nabla_{\theta} \pi_{\theta}(\tau) R(\tau) d\tau \\&= \mathbb{E}_{\tau \sim \pi_{\theta}} \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} R(\tau)\end{aligned}$$

# GRADIENT

Recall that:

$$\tau = (s_1, a_1, r_1, s_2, \dots, s_T, a_T, r_T, s_{\text{END}})$$

$$\pi_\theta(\tau) = p(s_1) \pi_\theta(a_1 | s_1) p(r_1, s_2 | s_1, a_1) \dots p(r_T, s_{\text{END}} | s_T, a_T)$$

Then:

$$\begin{aligned} \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} &= \frac{\nabla_\theta [p(s_1) \pi_\theta(a_1 | s_1) p(r_1, s_2 | s_1, a_1) \dots p(r_T, s_{\text{END}} | s_T, a_T)]}{p(s_1) \pi_\theta(a_1 | s_1) p(r_1, s_2 | s_1, a_1) \dots p(r_T, s_{\text{END}} | s_T, a_T)} \\ &= \frac{\nabla_\theta [\pi_\theta(a_1 | s_1) \dots \pi_\theta(a_T | s_T)]}{\pi_\theta(a_1 | s_1) \dots \pi_\theta(a_T | s_T)} \\ &= \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} \end{aligned}$$

## GRADIENT

We obtained:

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} R(\tau) = \mathbb{E}_{\tau \sim \pi_{\theta}} \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} R(\tau) = \mathbb{E}_{\tau \sim \pi_{\theta}} \sum_{t=1}^T \frac{\nabla_{\theta} \pi_{\theta}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} R(\tau)$$

# UNDERSTANDING THE GRADIENT OF THE EXPECTED RETURN

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

$\pi_\theta(\cdot | s_t)$  – output of the policy network on  $s_t$  (probability vector)

$\pi_\theta(a_t | s_t)$  – The probability of the action  $a_t$  that **actually happened** in time  $t$ .

$$\pi_\theta(a_t | s_t) = \pi_\theta(a^{(i)} | s_t) = \pi_\theta(\cdot | s_t)_i,$$

where  $a_t = a^{(i)}$  for some  $i \in \{1, \dots, A\}$ .

# UNDERSTANDING THE GRADIENT

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

For a given trajectory  $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T, s_{\text{END}})$ , we:

# UNDERSTANDING THE GRADIENT

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

For a given trajectory  $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T, s_{\text{END}})$ , we:

- Compute the gradient of the probability of the taken action  $a_t$  with respect to parameters  $\theta$ .

# UNDERSTANDING THE GRADIENT

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

For a given trajectory  $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T, s_{\text{END}})$ , we:

- Compute the gradient of the probability of the taken action  $a_t$  with respect to parameters  $\theta$ .
- Multiply it by  $R(\tau)$ .

# UNDERSTANDING THE GRADIENT

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

For a given trajectory  $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T, s_{\text{END}})$ , we:

- Compute the gradient of the probability of the taken action  $a_t$  with respect to parameters  $\theta$ .
- Multiply it by  $R(\tau)$ .
- Divide it by the value of the probability.

# UNDERSTANDING THE GRADIENT

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

For a given trajectory  $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T, s_{\text{END}})$ , we:

- Compute the gradient of the probability of the taken action  $a_t$  with respect to parameters  $\theta$ .
- Multiply it by  $R(\tau)$ .
- Divide it by the value of the probability.
- Repeat this for all  $t = 1, \dots, T$  and sum the results.

## GRADIENT – $\nabla_{\theta}\pi_{\theta}(a_t \mid s_t)$ INTUITION

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \sum_{t=1}^T \frac{\nabla_{\theta}\pi_{\theta}(a_t \mid s_t)}{\pi_{\theta}(a_t \mid s_t)} R(\tau)$$

- The agent in state  $s_t$  chose action  $a_t$ , by following a policy  $\pi_{\theta}$ .
- We want to update the parameters  $\theta$ , to obtain a better policy  $\pi_{\theta'}$ .
- The gradient  $\nabla_{\theta}\pi_{\theta}(a_t \mid s_t)$  describes the direction of the fastest increase of  $\pi_{\theta}(a_t \mid s_t)$ .
- Moving  $\theta$  in the direction of the gradient results in increasing the  $\pi_{\theta}(a_t \mid s_t)$  for the chosen action  $a_t$  and decreasing it for the other actions (as the probability sums to 1).

## GRADIENT – $R(\tau)$ INTUITION

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

Action  $a_t$  contributed to  $R(\tau)$ , so we want to change the probability of  $a_t$  *proportionally* to the return. In particular:

- When  $R(\tau) \geq 0$  we want to follow the gradient.
- When  $R(\tau) < 0$  we want to decrease  $\pi_\theta(a_t | s_t)$  – i.e. we want to move in direction opposite to the gradient.

$\nabla_\theta \pi_\theta(a_t | s_t) R(\tau)$  means:

"Modify  $\theta$  to change the probability of  $a_t$  proportionally to  $R(\tau)$ ."

## GRADIENT – DIVIDING BY $\pi_\theta(a_t | s_t)$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

Let  $\tau$  and  $\tau'$  be two different trajectories. Assume that:

$$R(\tau) = 2 \cdot R(\tau')$$

$$3 \cdot \pi_\theta(\tau) = \pi_\theta(\tau')$$

- The trajectory  $\tau$  is twice as good as  $\tau'$ , so obviously the agent should prefer  $\tau$  over  $\tau'$ .

## GRADIENT – DIVIDING BY $\pi_\theta(a_t | s_t)$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

Let  $\tau$  and  $\tau'$  be two different trajectories. Assume that:

$$R(\tau) = 2 \cdot R(\tau')$$
$$3 \cdot \pi_\theta(\tau) = \pi_\theta(\tau')$$

- The trajectory  $\tau$  is twice as good as  $\tau'$ , so obviously the agent should prefer  $\tau$  over  $\tau'$ .
- However  $\tau'$  happens 3 times more often than  $\tau$  in the training set, so it has the "effective weight" of  $3 \cdot R(\tau')$ .

## GRADIENT – DIVIDING BY $\pi_\theta(a_t | s_t)$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

Let  $\tau$  and  $\tau'$  be two different trajectories. Assume that:

$$R(\tau) = 2 \cdot R(\tau')$$
$$3 \cdot \pi_\theta(\tau) = \pi_\theta(\tau')$$

- The trajectory  $\tau$  is twice as good as  $\tau'$ , so obviously the agent should prefer  $\tau$  over  $\tau'$ .
- However  $\tau'$  happens 3 times more often than  $\tau$  in the training set, so it has the "effective weight" of  $3 \cdot R(\tau')$ .
- Without the normalization by the probability of the trajectories, the agent would learn to prefer  $\tau'$  over  $\tau$ !

## GRADIENT – SUMMING OVER $t$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau)$$

- Take sum over all  $t = 1, \dots, T$  to modify all preferences.
- Note that we do not take the average (i.e. divide by  $T$ ), because we want trajectories with high  $R(\tau)$ , even if they are long.

## LOGARITHM OF THE GRADIENT

Sometimes it is more desired (for numerical stability) to compute the logarithms (or logits) of the actions probabilities.

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} R(\tau) = \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R(\tau)$$

Note that in order to sample the actions we still need to compute the probabilities  $\pi_\theta(\cdot | s_t)$ .

## WEIGHTED GRADIENTS

Until now the gradient of our objective was described by:

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R(\tau)$$

However, we would like to work with:

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t$$

The difference is that in the second formula we weight the gradients  $\nabla_\theta \ln \pi_\theta(a_t | s_t)$  only with the rewards obtained **after** the action. We shall later argue that this approach improves the training.

ARE WE ALLOWED TO REPLACE  $R(\tau)$  WITH  $R_t$ ?

## WEIGHTED GRADIENTS – $R(\tau)$

Let's take a closer look at the gradient with  $R(\tau)$ :

$$\begin{aligned} & \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R(\tau) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left( \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \right) \cdot \left( \sum_{t=1}^T r_t \right) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t,t'=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'}. \end{aligned}$$

Observe that each gradient  $\nabla_\theta \ln \pi_\theta(a_t | s_t)$  is weighted by **every** reward of the trajectory.

## WEIGHTED GRADIENTS – $R_t$

Now let's investigate the gradient with  $R_t$ :

$$\begin{aligned} & \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left( \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \right) \cdot \left( \sum_{t'=t}^T r_{t'} \right) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \sum_{1 \leq t \leq t' \leq T} \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} \end{aligned}$$

Observe that each gradient  $\nabla_\theta \ln \pi_\theta(a_t | s_t)$  is weighted only by **future** rewards.

## WEIGHTED GRADIENTS – $R_t$

$$\begin{aligned}\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t,t'=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} &\stackrel{?}{=} \mathbb{E}_{\tau \sim \pi_\theta} \sum_{1 \leq t \leq t' \leq T} \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} \\ \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t>t'} \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} &\stackrel{?}{=} 0\end{aligned}$$

The answer is YES. It may be proved that if  $t > t'$ , then  $r_{t'}, a_t$  are independent given  $s_t$ , and:

$$\mathbb{E}_{\tau \sim \pi_\theta} \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} = 0.$$

## WEIGHTED GRADIENTS – $R_t$

We get that:

$$\mathbb{E}_{\tau \sim \pi_\theta} \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot r_{t'} = 0, \text{ for } t > t'$$

Why is directly excluding the above term from the entire gradient computation helpful?

- The above equality holds on average, but (usually) not on the training set.
- Therefore it only introduces noise.
- Also, leaving those terms is equal to giving  $R(\tau)$  as the last reward  $r_T$  (with all previous rewards set to 0).
- In consequence, we lose information on when the rewards were obtained.

## PRACTICAL NOTES – DISCOUNT FACTOR

Choose  $0 < \gamma < 1$  (usually  $0.9 < \gamma < 1$ ) and replace

$$R_t = r_t + r_{t+1} + r_{t+2} + \dots + r_T = \sum_{j=t}^T r_j$$

with

$$R_t^\gamma = \underline{\gamma^0} r_t + \underline{\gamma^1} r_{t+1} + \underline{\gamma^2} r_{t+2} + \dots + \underline{\gamma^{T-t}} r_T = \sum_{j=t}^T \gamma^{j-t} r_j.$$

Note that:

$$1 = \gamma^0 > \gamma^1 > \gamma^2 > \dots > \gamma^{T-t}$$

## DISCOUNT FACTOR – INTUITION

$$R_t^\gamma = \underline{\gamma^0 r_t} + \underline{\gamma^1 r_{t+1}} + \underline{\gamma^2 r_{t+2}} + \dots + \underline{\gamma^{T-t} r_T}$$

- If  $\gamma \rightarrow 1$ , then  $R_t^\gamma \rightarrow R_t$ 
  - We consider all rewards received after time  $t$ .
- If  $\gamma \rightarrow 0$ , then  $R_t^\gamma \rightarrow r_t$ 
  - We consider only the imminent reward.
- Rewards might be more relevant to recent actions.
  - $R_t^\gamma$  is a compromise.
  - Early rewards – weighted with  $\approx 1$ .
  - Late rewards – also important, but added with a smaller weight.
- We reduce noise, but introduce error (bias).
  - Still, it's beneficial.

## DISCOUNT FACTOR – HORIZON

We define the horizon  $h_\gamma$  as:

$$h_\gamma := \frac{1}{1 - \gamma}.$$

The horizon indicates how far ahead we look:

- For  $\gamma \approx 1$  we have

$$\gamma^{h_\gamma} \approx \frac{1}{e}, \quad \gamma^{3h_\gamma} \approx 0.05$$

- Examples

$$h_{0.8} = 5, \quad h_{0.95} = 20, \quad h_{0.999} = 1000$$

- In practice: rather than choosing  $\gamma$  first decide for the horizon  $h$  and then set

$$\gamma = 1 - \frac{1}{h}$$

## PRACTICAL NOTES – OBJECTIVE FUNCTION $J(\theta)$

Let  $\{\tau_1, \dots, \tau_N\}$  is the training set of  $N$  trajectories generated by policy  $\pi_\theta$ . The approximation of the gradient is:

$$\begin{aligned}\nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} R(\tau) &= \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t \\ &\approx \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t^\gamma \\ &\approx \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t(\tau_j) | s_t(\tau_j)) \cdot R_t^\gamma(\tau_j)\end{aligned}$$

## PRACTICAL NOTES – OBJECTIVE FUNCTION $J(\theta)$

$$\frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \nabla_{\theta} \ln \pi_{\theta}(a_t(\tau_j) | s_t(\tau_j)) \cdot R_t^{\gamma}(\tau_j)$$

- Inefficient in practice – we have to separately calculate multiple gradients  $\nabla_{\theta} \ln \pi_{\theta}(a_t | s_t)$ .
- PyTorch and TensorFlow do not like that.
- Observe that we can move the gradient symbol:

$$\nabla_{\theta} \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \ln \pi_{\theta}(a_t(\tau_j) | s_t(\tau_j)) \cdot R_t^{\gamma}(\tau_j) =: \nabla_{\theta} J(\theta)$$

- Then we only have to calculate  $J(\theta)$  and ask PyTorch to compute the gradient and perform the ascent step.

## TRAINING LOOP

1. Sample  $N$  trajectories  $\tau_1, \dots, \tau_N$ .
2. Calculate the function

$$-J(\theta) = -\frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \ln \pi_\theta(a_t(\tau_j) | s_t(\tau_j)) \cdot R_t^\gamma(\tau_j).$$

3. Ask PyTorch to minimize this function by performing gradient descent step.

## PART 3: VALUE FUNCTION

# PG WITH BASELINE

## ENVIRONMENT $E$

Let  $E$  be an environment, with discrete actions

$\mathbb{A} = \{a^{(1)}, a^{(2)}, a^{(3)}\}$  and states  $s \in \mathbb{S}$ . Let  $s_t$  be the state visited at time  $t$  and  $R_t$  the corresponding return.

## ENVIRONMENT $E'$

Now consider an environment  $E'$ , which is basically a copy of  $E$ .

The only difference is that, for every trajectory, after reaching the ending state  $s_{\text{END}}$  we subtract 10 from the trajectories return.

## PG WITH BASELINE

Let  $R_t$ ,  $r_t$ ,  $R'_t$ ,  $r'_t$  be the returns and rewards for state  $s_t$  in environments  $E$  and  $E'$ , respectively. Note that:

- The rewards stay the same:  $r'_t = r_t$ .
- The **returns** change:

$$R'_t = R_t - 10.$$

- The optimal policy  $\pi_{optim}$  in  $E$  is also an optimal policy in  $E'$  and *vice versa*.

## PG WITH BASELINE

- On **average** there is no difference in how the agent learns in those two environments:

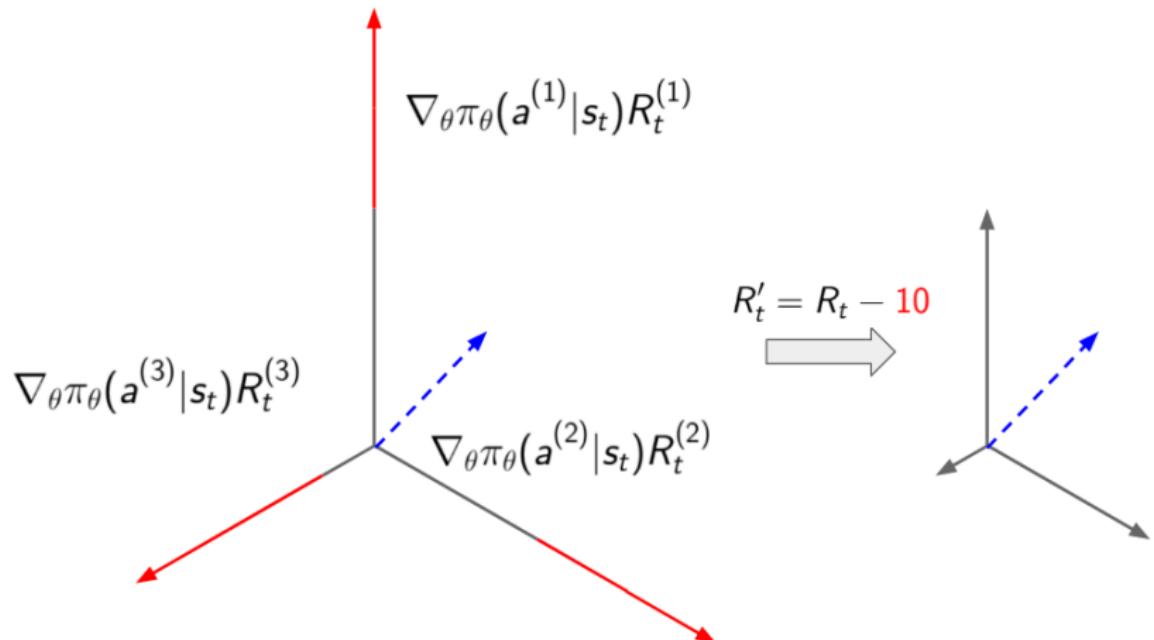
$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} R'_t \right] = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} (R_t - 10) \right],$$

since:

$$\begin{aligned} \mathbb{E}_{\tau \sim \pi_\theta} \left[ \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} 10 \right] &= 10 \cdot \mathbb{E}_{\tau \sim \pi_\theta} \left[ \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} \right] = \\ &10 \cdot \int \nabla_\theta \pi_\theta(\tau) d\tau = 10 \nabla_\theta 1 = 0 \end{aligned}$$

- In **practice** however, learning in  $E'$  may be more stable.

## PG WITH BASELINE



## BASELINE FUNCTION

It turns out that we may subtract from the returns any **baseline** function  $b : \mathbb{S} \rightarrow \mathbb{R}$

$$\begin{aligned}\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t &= \\ \mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot (R_t - b(s_t))\end{aligned}$$

It may be proved that:

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot b(s_t) = 0$$

The proof is similar to the simple constant case  $b(s_t) = 10$ .

## BASELINE FUNCTION

Note that we are changing only the **returns**, **NOT** the rewards.

- Increasing all rewards for a given state  $s_t$  would result in agent visiting this state more often. In consequence we would change the objective – BAD IDEA!
- Modifying the return is a trick that always works.
- Note that we can select a different baseline value for each state. For example we can decrease the returns by 10 in one state and increase by 20 in another.
- Unbiased noise reduction – recall replacing  $R(\tau)$  with  $R_t$ .

### HOW TO DEFINE A GOOD BASELINE?

## VALUE FUNCTION

The value function  $V_{\pi_\theta} : \mathbb{S} \rightarrow \mathbb{R}$  is defined as:

$$V_{\pi_\theta}(s) := \mathbb{E}_{\tau \sim \pi_\theta, s_1(\tau)=s} R(\tau),$$

where  $s_1(\tau)$  denotes the first state of  $\tau$ .

The value function depends on *policy* in a complicated way:

- It also involves the environment.
- It cannot be calculated.

# VALUE FUNCTION – MARKOV PROPERTY!

The value of a state does not change in time:

$$\mathbb{E}_{\tau \sim \pi_\theta, s_1(\tau) = s} R(\tau) = \mathbb{E}_{\tau \sim \pi_\theta, s_t(\tau) = s} R_t(\tau)$$

Intuition:

- Agent currently is in state  $s$ , time  $t$  is irrelevant.
- How much rewards in total do we expect from now on?

We can learn the value function from sampled trajectories.

## VALUE NETWORK

We will approximate the true value function  $V_{\pi_\theta}$  with a neural network  $V_\psi$  with weights  $\psi$ :

$$V_\psi \approx V_{\pi_\theta}.$$

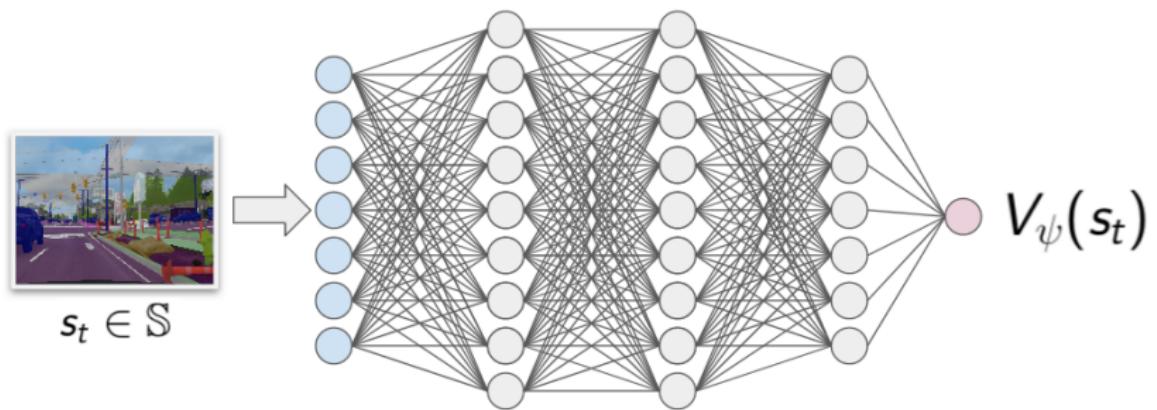
The training set consists of pairs  $(s_t, R_t)$ :

- May be formed from a single trajectory:

$$\{(s_1, R_1), \dots, (s_T, R_T)\}$$

- We can also combine multiple trajectories.

# VALUE NETWORK



## VALUE NETWORK: OBJECTIVE

Learning the value network is a regression task. The cost function of the value network is the standard *Mean Squared Error*.

Recall that the value function depends on the policy! Therefore the value network must be tuned after every update of  $\pi_\theta$ :

- Different  $\pi_\theta$  means different target  $V_{\pi_\theta}$ .
- In contrast to Policy Gradient, we can do multiple updates of  $\psi$  without resampling the training set.

## ADVANTAGE LEARNING

Let's define the advantage at time  $t$  as:

$$A_t := R_t - V_\psi(s_t)$$

Replacing the return  $R_t$  in the gradient estimation with  $A_t$  is equivalent to using  $V_\psi$  as a **baseline**:

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot R_t =$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot (R_t - V_\psi(s_t)) =$$

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot A_t$$

$$A_t := R_t - V_\psi(s_t)$$

Note that:

- $R_t$  is the absolute gain.  $A_t$  is the gain relative to the average.
- The advantage gives a better estimator for finite number of trajectories.
- Remember – any baseline is correct!
  - $V_\psi$  doesn't even have to be a good approximation of  $V_{\pi_\theta}$ .
  - However better approximation gives more stability!

## TRAINING LOOP – POLICY GRADIENT WITH BASELINE

1. Sample  $N$  trajectories  $\tau_1, \dots, \tau_N$ .
2. Improve the parameters of our value network by minimizing the loss (apply multiple gradient descent steps):

$$\frac{1}{M} \sum_{j=1}^N \sum_{t=1}^T (R_t(\tau_j) - V_\psi(s_t(\tau_j)))^2,$$

where  $M$  is the total number of samples in all trajectories.

3. Calculate the function

$$J(\theta) = \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \ln \pi_\theta(a_t(\tau_j) | s_t(\tau_j)) \cdot [R_t(\tau_j) - V_\psi(s_t(\tau_j))],$$

4. Ask PyTorch to minimize the function  $-J(\theta)$  by performing gradient descent step.

## PART 4: PROXIMAL POLICY OPTIMIZATION

# REMARKS ON SOME COMMON PROBLEMS

## SAMPLING TRAJECTORIES

Collecting new trajectories is the most expensive part of RL!

- We would like to reuse trajectories for multiple updates of the agent.
- Recall that this is allowed for  $V_\psi$ , but not for  $\pi_\theta$ .

## DISTRIBUTION OF TRAJECTORIES CHANGES DURING THE LEARNING

The update of the agent changes the distribution of the trajectories. For example:

- Assume  $(s_t, a_t)$  has a positive return. Then after the update  $\pi_\theta(a_t | s_t)$  may increase.
- This means that the pairs  $(s_t, a_t)$  are now **undersampled** in the training set!

## QUESTION

IS IT POSSIBLE TO MAKE MULTIPLE UPDATES OF  $\pi_\theta$  WITHOUT RESAMPLING TRAJECTORIES?

# PROXIMAL POLICY OPTIMIZATION

Sketch of the algorithm:

1. Copy  $\pi_\theta$  to  $\pi_{\theta\text{OLD}}$ .
2. Sample several trajectories using  $\pi_{\theta\text{OLD}}$ .
3. Make *multiple updates* of  $\pi_\theta$ .

**PROBLEM:** Until now we assumed  $\pi_\theta = \pi_{\theta\text{OLD}}$  and only *one update* of  $\pi_\theta$  per training set. Therefore the equations need to be modified!

# PROXIMAL POLICY OPTIMIZATION

Probability that  $(s_t, a_t)$  will occur in the training set is a product of:

- $\pi_\theta(a_t | s_t)$  – probability of taking action  $a_t$  after observing  $s_t$ .
  - Explicitly given by the policy network.
- $\pi_\theta(s_t)$  – probability of visiting  $s_t$ .
  - Depends on the environment and  $\pi_\theta(a_{t'} | s_{t'})$ , where  $t' < t$ .
  - Different actions *lead to* different states later on.

# PROXIMAL POLICY OPTIMIZATION

## THE PROXIMAL POLICY OPTIMIZATION ASSUMPTION

As long as:

$$1 - \epsilon < \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta\text{OLD}}(a_t | s_t)} < 1 + \epsilon,$$

for some (small, but significant)  $\epsilon$ , e.g.  $\epsilon = 0.2$ , we may assume that  $\pi_\theta(s_t)$  is close to  $\pi_{\theta\text{OLD}}(s_t)$ :

$$\pi_\theta(s_t) \approx \pi_{\theta\text{OLD}}(s_t)$$

Therefore, the following is a good approximation:

$$\mathbb{E}_{\tau \sim \pi_\theta} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} \approx \mathbb{E}_{\tau \sim \pi_{\theta\text{OLD}}} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_{\theta\text{OLD}}(a_t | s_t)}$$

## PPO ASSUMPTION – SOME INTUITION

$$\mathbb{E}_{\tau \sim \pi_\theta} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} \approx \mathbb{E}_{\tau \sim \pi_{\theta \text{OLD}}} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_{\theta \text{OLD}}(a_t | s_t)}$$

- The left part is what we want.
- However, we must use  $\mathbb{E}_{\tau \sim \pi_{\theta \text{OLD}}}$ , because training set was sampled from  $\pi_{\theta \text{OLD}}$ .
- Recall that we always divide by the probability of being in the training set. Therefore we also substitute the denominator. This is an important step, since difference between  $\pi_\theta(a_t | s_t)$  and  $\pi_{\theta \text{OLD}}(a_t | s_t)$  is relevant ( $\epsilon \not\approx 0$ ).

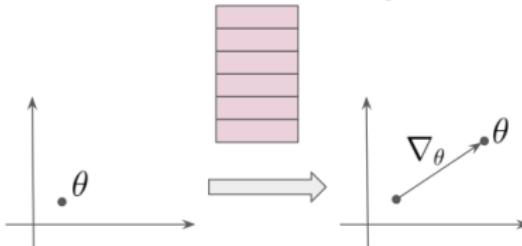
## PPO ASSUMPTION – SOME INTUITION

$$\mathbb{E}_{\tau \sim \pi_\theta} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} \approx \mathbb{E}_{\tau \sim \pi_{\theta \text{OLD}}} \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_{\theta \text{OLD}}(a_t | s_t)}$$

- We keep the numerator, because we update  $\pi_\theta$ , not  $\pi_{\theta \text{OLD}}$ .
- We do not have equality because  $\pi_\theta(s_t) \neq \pi_{\theta \text{OLD}}(s_t)$ .  
However we assume that this part of error is *relatively* small.
  - We could introduce the necessary corrections, but then the estimator would be very unstable.
  - Analogically to the discounting trick – we introduce error, but gain stability!

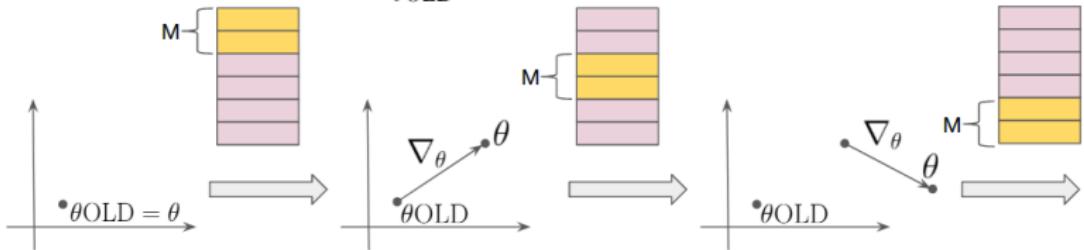
## Policy Gradient

SAMPLED DATASET  $\sim \pi_\theta$



## Proximal Policy Optimization

SAMPLED DATASET  $\sim \pi_{\theta \text{OLD}}$



# PROXIMAL POLICY OPTIMIZATION

Let's define  $\rho_t$  as:

$$\rho_t := \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta\text{OLD}}(a_t | s_t)}.$$

This simplifies the notation (denominator is constant wrt. to  $\theta$ ):

$$\nabla_\theta \rho_t = \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_{\theta\text{OLD}}(a_t | s_t)}.$$

For every  $t$  we try to keep

$$1 - \epsilon < \rho_t < 1 + \epsilon,$$

so that we can use the formula:

$$\mathbb{E}_{\tau \sim \pi_\theta} \sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_t | s_t) \cdot A_t \approx \mathbb{E}_{\tau \sim \pi_{\theta\text{OLD}}} \sum_{t=1}^T \nabla_\theta \rho_t \cdot A_t$$

## WHY $\rho_t \neq 1$ ?

$$\rho_t = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta\text{OLD}}(a_t | s_t)}$$

- initially  $\pi_\theta(a_t | s_t) = \pi_{\theta\text{OLD}}(a_t | s_t)$ , so  $\rho_t = 1$ .
- $\pi_{\theta\text{OLD}}(a_t | s_t)$  is constant, thus  $\rho_t$  changes with  $\pi_\theta(a_t | s_t)$ .

$$\mathbb{E}_{\tau \sim \pi_{\theta\text{OLD}}} \sum_{t=1}^T \nabla_\theta \rho_t \cdot A_t$$

- if  $A_t > 0$ , then  $\rho_t$  will usually increase (decrease with  $A_t < 0$ )
  - Every time it's included in training batch, gradient will try to increase  $\pi_\theta(a_t | s_t)$ .
  - Caution! other training samples will have unknown influence on this particular  $\rho_t$ .

## SURROGATE OBJECTIVE

For convenience we define a surrogate objective

$$\mathcal{J} := \mathbb{E}_{\tau \sim \pi_{\theta}^{\text{OLD}}} \sum_{t=1}^T \rho_t A_t,$$

and we have

$$\nabla_{\theta} \mathcal{J} = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}^{\text{OLD}}} \sum_{t=1}^T \rho_t \cdot A_t = \mathbb{E}_{\tau \sim \pi_{\theta}^{\text{OLD}}} \sum_{t=1}^T \nabla_{\theta} \rho_t \cdot A_t.$$

Note that we can simply swap  $\nabla_{\theta}$  and  $\mathbb{E}_{\tau \sim \pi_{\theta}^{\text{OLD}}}$ , because  $\pi_{\theta}^{\text{OLD}}$  does **not** depend on  $\theta$ .

## HOW TO ENFORCE $1 - \epsilon < \rho_t < 1 + \epsilon$ ?

- We always keep a copy of  $\pi_{\theta^{\text{OLD}}}(a_t | s_t)$ , so  $\rho_t$  may be calculated if needed.
- If  $\rho_t$  goes out of bounds, we temporarily remove the corresponding training example from the training set.
- Important change: if  $A_t > 0$ , then we only enforce  $\rho_t < 1 + \epsilon$ 
  - We *do not* want to exclude examples that went out of range by accident.
  - Analogically for  $A_t < 0$  enforce only  $\rho_t > 1 - \epsilon$ .
- Turns out if we modify our surrogate objective a bit, then we can achieve what we want:

$$\tilde{\mathcal{J}} := \mathbb{E}_{\tau \sim \pi_{\theta^{\text{OLD}}}} \sum_{t=1}^T \left[ \min(\rho_t \cdot A_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \cdot A_t) \right]$$

## THE SURROGATE OBJECTIVE – INTUITION

$$\min (\rho_t \cdot A_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \cdot A_t)$$

Additional intuition taken from the original PPO paper:

- The first part,  $\rho_t \cdot A_t$ , comes from  $\mathcal{J}$ .
- "The second term (...) removes the incentive of moving  $\rho_t$  outside of the interval  $[1 - \epsilon, 1 + \epsilon]$ ."
- "We take the minimum of the clipped and unclipped objective, so the final objective is a lower bound (i.e., a pessimistic bound) on the unclipped objective."

## EXERCISE

Prove (on paper) that this formula works as intended

$$\min(\rho_t \cdot A_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \cdot A_t).$$

Check all 6 possible cases:

- $A_t$  may be positive or negative.
- $\rho_t$  may be  $< 1 - \epsilon$ ,  $> 1 + \epsilon$ , or  $\approx 1$ .

Hints:

- Make two plots of the formula wrt. to  $\rho_t$  for positive and negative  $A_t$ .
- If  $\rho_t$  is out of bounds, then gradient of the second part is equal to zero.

## TRAINING LOOP – PPO

1. Sample  $N$  trajectories  $\tau_1, \dots, \tau_N$  from  $\pi_\theta$ .
2. Assign  $\pi_{\theta\text{OLD}} := \pi_\theta$ .
3. Loop over batches  $\mathcal{B}$  from gathered trajectories:
  - 3.1 Improve the parameters of our value network by minimizing the loss:

$$\frac{1}{|\mathcal{B}|} \sum_{(s_i, R_i) \in \mathcal{B}} [R_i - V_\psi(s_i)]^2.$$

- 3.2 Calculate the function:

$$\tilde{J}(\theta) = \frac{1}{|\mathcal{B}|} \sum_{(s_i, a_i, R_i) \in \mathcal{B}} \left[ \min \left( \rho_i \cdot A_i, \text{clip}(\rho_i, 1 - \epsilon, 1 + \epsilon) \cdot A_i \right) \right].$$

- 3.3 Ask PyTorch to minimize the function  $-\tilde{J}(\theta)$  by performing gradient descent step.

## FURTHER READINGS

1. "Reinforcement Learning: An Introduction"  
(Sutton, Barto)
2. UCL Course on RL  
(Silver, [videos](#))
3. Deep Reinforcement Learning course at UC Berkeley  
(Levine, [videos](#))
4. Proximal Policy Optimization  
(Schulman, Wolski, Dhariwal, Radford, Klimov)
5. OpenAI Resources page
6. The Policy of Truth  
(Recht, "policy gradient is nothing more than random search dressed up in mathematical symbols and lingo")