

Ruggero Rossi
r.rossi@opencomplexity.com

A Gentle Principled Introduction to Deep Reinforcement Learning

Contents

0. Introduction.....	2
1. Reinforcement Learning Definition.....	2
2. Reinforcement Learning Formalization	6
3. Variety of Methods.....	9
4. Q-Learning	11
5. Policy Based Methods	21
6. Policy Gradient	22
7. Generalized Advantage Estimation.....	34
8. Natural Policy Gradient.....	36
9. Trust Region Policy Optimization (incomplete)	41
10. Proximal Policy Optimization	49
11. Reward Shaping and Curriculum Learning	53
12. Imitation Learning	54
13. Afterword	57
Appendix A: Information Theory Refresh	58
References	67

0. Introduction

This paper is an attempt to introduce Deep Reinforcement Learning to those who possess a solid knowledge of Deep Learning. The aim is to explain the functioning of selected Deep Reinforcement Learning algorithms that are both practically effective and didactically representative, providing the reader with a theoretical justification and a balanced level of proofs and mathematical details, without losing sight of the general picture. The prerequisite is the knowledge of Deep Learning, that implies knowledge of basic Calculus, Linear Algebra, Probability and Statistics.

1. Reinforcement Learning Definition

Reinforcement Learning is a method for an artificial agent to do automated learning using the outcomes caused by its actions. A definition from [Sutton & Barto 2018] reports: “Reinforcement learning is learning what to do - how to map situations to actions - so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them”.

So, paraphrasing Sutton and Barto, Reinforcement Learning is: learning which action to take, depending on the current situation, with the purpose of maximizing the total reward.

In that definition we have some of the base components of a Reinforcement Learning system:

- an *agent* able to decide which action to take
- a *state* of the world – or better an *observation* about the state of the world
- a *reward*

The *agent* is obviously the artificial intelligence that is trying to learn how to act, it may be just a program inside a simulation as well as a physical robot in a real environment.

The *observation* about the state of the world is the data that the agent possesses about the state of world, and that may be potentially incomplete, noisy, or delayed. The observation may be given to the agent by the simulator software, in case of simulated environment, or may be obtained by the agent through sensors in case of physical robot.

The *reward* is a number that represents how well or how badly the agent is behaving (usually it is positive for good rewards and negative for bad rewards). It is something similar to the concept of utility. In case of a simulated environment the reward is given to the agent by the simulator itself, together with the timed information updates such as the observation of the state of the world. In case of physical robot in a real environment instead, a software layer must be programmed in the robot in order to understand what is happening in the world and calculate a reward that is related to how much good or bad the situation is.

In both cases the reward is a number whose value is communicated to the agent by a human-made program: either in the simulation or in the robot software. There is not such a thing in the environment as a “reward number”, the reward is just a fictitious number that permits to do numerical optimization (humans and animals feel pain and pleasure and that may be seen as a reward system that has been evolved, but for our artificial agents it has to be built by us). Anyway in some cases that number may be somehow directly available from the real world or from the final purpose objective: think about an artificial intelligence trying to maximize the gains in the stock market, the reward may be directly the market returns, or think about of a robot trying to pick up trash from the ground, a reward proportional to garbage weight may be given every time that garbage is actually collected.

If the task to be accomplished by the agent is never ending it is said to be “infinite-horizon”, otherwise if a begin and an end of the sequence of actions can be identified, it is said “finite-horizon” and we may call “episode” what happens within the begin and the end. For instance, if a robot is trying to shoot a basketball inside a hoop, an episode starts when the robot tries to grab the ball and ends when the ball, after being shot, lands, either passing through the hoop or outside it. An episode contains a “*trajectory*” made of a sequence of states and actions. The sum of the rewards obtained in an “*episode*”, possibly discounted by a time distance value, is called “*return*” (the time discount represents the preference of receiving rewards immediately rather than in the future).

Another useful concept when talking about reinforcement learning is the concept of “*Policy*”: with that term we mean a strategy that associates an action to a world state, for each possible state. That means that in each situation an agent is, the Policy will tell him what action to take. A good Policy is one that makes the agent collect a good amount of rewards, an optimal policy is one that makes the agent collect the maximum amount of rewards (there may be more than one optimal policy). A policy may be stochastic and define, for each state, a probabilistic mixture of actions to take, such as “take action A with probability 80%, take action B with probability 20%”.

So, what finally an agent wants to learn is an optimal policy (or at least one good enough ! If finding the optimal policy takes too much time or resources, it may be better to find just a decent policy and use that, that is the topic “exploration vs exploitation” discussed later).

The agent may or may not have knowledge of how the world works and what happens if he takes a certain action, in which state he will end up and what reward will be obtained because of that action. When the agent has that kind of knowledge it is possible to apply the so-called *model-based* reinforcement learning algorithms. The model of the world is generally considered stochastic, with deterministic models being just particular cases. The model may be a fixed, prior knowledge of the agent. If instead the agent has not a prior model of the world mechanics, it is possible either to learn one in the process, or to use *model-free* algorithms that don't require a model of the world. In some case the model of the world is available (at least partially) because the task to be executed by the agent is to solve a game, such as Chess or Go, and the rules and mechanics of the game are known and may be taken in account by the algorithm. In some other cases the learning algorithm may be trying to learn a model of the world through experience and then use that model to apply a model-based algorithm. When the model of the world is complete, and when there are a finite number of states it is theoretically possible to calculate the best action for each state just with numerical optimization (e.g. with dynamic programming), using the “Bellman Equation” [Sutton and Barto 2018] without the need to make the agent taking real actions. Usually that is not the case because often the model of the world is not known (and sometimes difficult to be learned) and there is an infinite number of states, hence agents must advance through trial and errors to discover what is the best action in which situation.

Trying actions in the environment with the purpose of learning exposes the agents to the risk of obtaining very bad rewards: in a real environment that would mean damaging seriously the robot, or even worse creating risky situations outside the experimental environment. So, there is a trade-off between exploration (trying new actions to figure out if they bring better rewards) and exploitation (doing only the actions that so far showed to be sufficiently rewarding). Exploration exposes to risks but permits to optimize the policy, exploitation permits to gain the fruits of past exploration but avoids further learning. Agents usually are given a greater degree of exploration at the beginning, and the exploration degree is decreased as the learning progresses, favouring exploitation.

Some Reinforcement Learning algorithms (but not all) use the concept of “Value Function”: a function that takes the current state (or current observation) as input and evaluates the hypothetical goodness of it.

The evaluation score returned by the Value Function is related to how good it is expected to be the future situation in the long run, not just in the immediate next moments. If a certain state is expected to make the agent obtaining a reward, the value function will take in account the reward expected in that situation and all the rewards expected from what hypothetically may be the future situations if the agent acts as its “policy” suggests, possibly discounting the future rewards by a time-dependent distance value (rewards further in time may weight less than immediate rewards).

So, as I will describe later formally with the Bellman Equation, the Value function is a recursive concept: saying that the value at some state depends not only on the reward obtained in that state but also on all the rewards that may be obtained afterwards, is the same of saying that the value at a certain state depends not only on the goodness of that state “per se” but also on the values of the states that are reachable from there. This also implies that the value function depends on the policy: two different policies may reach different states from the same starting point and hence they may have two different value functions. So, in general each value function is tied to a certain policy: changing the policy changes the value function.

There are two different types of value functions: proper “value functions” and “action-value functions”.

The proper “value function”, returns the value of a state, considering that the agent will follow the policy from that state, so its only input is the state (and implicitly the policy). It is sometimes called also “state-value” function.

The “action-value” function takes as input the state and also an action, and returns a value considering that the agent will take that particular action in that state (even if it’s a different action with respect to what the policy would suggest), and after that it will follow his policy.

It has to be noted that when in literature it is used the term “Value Function” it may be referred to two very different things: the first is the value function as described above, the second is the estimate of the value function that an agent is approximating, and that may be far from the real value (and it should be better referred to with the terms “value function estimate”). When there is a finite number of states the estimate of the value function is often saved in a table with an estimate for each state (or, in case of “action-value”, an estimate for each state-action pair), when instead there is a great number of states that cannot be done and a function approximator is used, such as a neural network. That may be the case when the state/observation is described by continuous variables which if discretized in a table would need too much memory, or which would lose relevant details.

The term “Deep Reinforcement Learning” is used to describe Reinforcement Learning methods that use neural networks, especially neural networks with many hidden layers (“deep”). For instance, neural networks may be used to approximate value function estimates, or, as we will see later, for policy functions. The usage of deep neural networks allows to learn complex policies that are not possible with linear methods as function approximators. Moreover, they allow to learn end-to-end, that means having a neural network which, without the hand-made engineering of different other software modules, receives the sensory input and internally learns all the necessary functionalities to process it and take the output action, as opposed to having a separate module that analyses the input, a next module that plans what to do, a further module that learns how to actuate the robot engines, etc.

2. Reinforcement Learning Formalization

A formalization of the Reinforcement Learning problem may be done expressing it as a Markov Decision Process (MDP). A MDP is a way to describe how an agent passes from one state to another and possibly obtains a reward as a consequence of its actions. One requirement is that states must have the “Markov Property”: the probabilities to move from one state to another and to obtain rewards depend only on the knowledge of the current state and the chosen action, and not on any previously accessed state. In other words, a state description has the Markov property if it contains all the necessary information from the past and from the present to determine the transition probabilities to other states and the rewards.

This means that theoretically, non-Markov states can be turned in Markov states if all the history of past states and interactions are added to them (but this way to describe states as long sequences of the past are not simple to be managed by algorithms).

In a MDP, the mechanism that makes an agents passing from a state to another as a consequence of its actions is described by a “*transition function*” that is generally probabilistic and it is not necessarily known by the agent (the agent may just experience the change in state after an action). The mechanism that assigns a reward depending on a certain state, an action executed in that state, and a consequent state, is described by a “*reward function*”, that is not necessarily known by the agent (the agent may just notice the obtained reward after every action).

The theoretical exposition of RL that follows, when a different source is not referenced, is taken from [Sutton & Barto 2018] and [OpenAi 2018A], with some change of variables names to have a consistent description.

Formally, following OpenAI definition, a MDP is a 5-tuple $\langle S, A, R, P, \rho_0 \rangle$:

- S is the set of all states
- A is the set of actions
- $R : S \times A \times S \rightarrow \mathbb{R}$ is the reward function, with $r_t = R(s_t, a_t, s_{t+1})$
- $P : S \times A \rightarrow \mathcal{P}(S)$ is the transition function, with $P(s'|s, a)$ being the probability of transitioning from state s to state s' after taking action a
- ρ_0 is the distribution of initial state

The policy that decides which action to take depending on the state, is a stochastic function π , such that if at time t the state is s_t , it will return a probability distribution over the actions, from which it may be sampled the action to take a_t :

$$a_t \sim \pi(\cdot | s_t)$$

Equation 2.1

A sequence of states and actions is called “trajectory”, identified by the symbol τ

$$\tau = (s_0, a_0, s_1, a_1, \dots)$$

Equation 2.2

The return $G(\tau)$ of a trajectory is the (potentially infinite) sum of all rewards of the trajectory, with $\gamma \in [0,1]$, and it is time-discounted if $\gamma < 1$:

$$G(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$$

Equation 2.3

Given a policy π , the probability of following any trajectory τ made of T steps is:

$$P(\tau|\pi) = \rho_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$$

Equation 2.4

And the expected return $J(\pi)$ is :

$$J(\pi) = \int_{\tau} P(\tau|\pi)G(\tau) = E_{\tau \sim \pi}[G(\tau)]$$

Equation 2.5

Where a notation like $E_{x \sim z}[f(x)]$ means: Expected Value of $f(x)$ where x follows the probability density function (or probability mass function) $z(x)$.

Later we will see also a notation like $E_{x \sim z}[f(x)|y = m, u = n]$ that means: Expected Value of $f(x)$ where x follows the pdf/pmf $z(x)$, given that y is equal to m and u is equal to n .

The problem of Reinforcement Learning is to find the policy that maximizes that expectation in eq. 2.5. Hence:

$$\pi^* = \arg \max_{\pi} J(\pi)$$

Equation 2.6

Where π^* is the optimal policy.

The Value Function for a given policy π is:

$$V^{\pi}(s) = E_{\tau \sim \pi}[G(\tau)|s_0 = s]$$

Equation 2.7

The Action-Value function (also named “Q-Value”) is:

$$Q^{\pi}(s, a) = E_{\tau \sim \pi}[G(\tau)|s_0 = s, a_0 = a]$$

Equation 2.8

It has to be noted that

$$V^{\pi}(s) = E_{a \sim \pi}[Q^{\pi}(s, a)]$$

Equation 2.9

As previously anticipated, the value function is a recursive concept because the value of a state is dependent on the value of the next state, recursively. The Bellman Equations express this relationship:

$$V^\pi(s) = E_{a \sim \pi} [R(s, a, s') + \gamma V^\pi(s')]_{s' \sim P}$$

Equation 2.10

$$Q^\pi(s, a) = E_{s' \sim P} [R(s, a, s') + \gamma E_{a' \sim \pi} [Q^\pi(s', a')]]$$

Equation 2.11

Using eq. 2.9 in 2.11 we obtain:

$$Q^\pi(s, a) = E_{s' \sim P} [R(s, a, s') + \gamma V^\pi(s')]$$

Equation 2.12

Some algorithms use the “Advantage Function” [Baird 1993], that is a function indicating how much it is better or worse to take a certain action (and after that follow the policy) instead of taking the action suggested by the policy, i.e. it indicates the relative advantage of taking a certain action with respect to the policy. The advantage function is the following:

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

Equation 2.13

3. Variety of Methods

Among RL methods, a first distinction can be made between **model based** and **model free** methods. Model based methods use a model of the world, consisting generally in the knowledge of the reward and transition functions, and the initial state distribution: R, P, ρ_0 . Model free methods instead do not use a model (generally in real environments a model is not available). Some other methods try to build a model from the observations. In this guide I will mainly deal with model free methods.

A second distinction is between on-policy and off-policy methods: **on-policy** methods evaluate and improve the same policy that is followed by the agent during training, while in **off-policy** methods the training is done following a policy that is different from the one that is evaluated or improved. In off-policy methods the policy that is being learned is named “*target policy*” and the policy that is used to generate samples is named “*behavior policy*”.

A third distinction may be made between value based methods and policy based methods. In **value based** methods the algorithm aims at estimating a value function or an action-value function and consequently modify the policy to prefer actions that lead to better values. The “*Policy Improvement Theorem*” [Sutton & Barto 2018, Ch.4] guarantees that modifying the policy for a state in a way that obtains a value improvement in that state, strictly improves the expected return $J(\pi)$ and hence improves the overall policy.

If a world model is not available, “proper” value functions (i.e. state-only values) do not contain enough information to improve the policy, hence the only model-free value based methods are those that use action-values: so they can choose the action that is associated with the best outcome for each state.

In **policy based** methods instead, the action is directly chosen by a policy function (that outputs a probability distribution over actions) and value or action-value functions may not be present, or may be used only to improve learning.

A fourth distinction of RL methods can be done between “tabular” and “approximate” methods. **Tabular** methods are those whose state space is composed by a finite set of discrete states, that are hence representable in memory by tables, for instance the Value Function estimation may be simply a table that associates a state to a value. **Approximate** methods are those whose state space is continuous or too big, and so its representation in the algorithm is made by some function approximation. For instance, the Value Function could be approximated (or estimated) by a neural network that takes the information about the state/observation (e.g. sensor data, webcam frame image, etc.) as input, and computes an approximate value as output. It is common to use convolutional neural networks as approximators when the state/observation input is an image. While tabular and approximate methods share much theory and algorithms, they also have important differences, the most notable of which is the fact that an approximate value based algorithm, if it is also using bootstrapping (i.e. using the approximate value function estimation as a replacement for the true value of value function in the Bellman equation, or equivalently using the approximated q-value estimation instead of true q-value) and if it is off-policy, it is not guaranteed to be convergent, while a tabular method would be, see [Sutton and Barto 2018, Ch. 11].

When a deep neural network is used as a function approximator in a RL algorithm, it is said to be Deep Reinforcement Learning.

4. Q-Learning

An example of (action-)value based method is the Q-Learning [Watkins 1989]. Let us see how it works.

The basic idea of Q-Learning is to estimate the value of any action-state pair, and then the target policy would be the policy that, given a state, selects the action with greater action-state value estimate. But the important detail in all this is that during learning, Q-Learning “bootstraps”, that means that it estimates the value of a past action summing the obtained reward for the past action with the estimate of the value function of the next state. This will be explained better later.

Q-Learning is an off-policy method, that means that the behavior policy (the policy used to generate trajectories) does not have to coincide with the target policy (the policy that we want to improve until it is optimal). In fact they could be very different, but for practical reasons it is better that they are not too much: since a function approximator is used to evaluate state-action value estimates, it is better to train that approximator on state-action samples coming from a state-action distribution that is similar to the one of the target policy. This is because the approximation capacity of a function approximator is limited and it is more precise when it is run on inputs similar to the ones used for training.

For this reason in Q-Learning the policy followed by agents during training (the *behavior policy*) is often made not too different from the target policy: the behavior policy sometimes chooses the action with the greatest action-value (exploitation) and sometimes a random action (exploration). A way to do that is to have an ϵ -greedy policy: with probability ϵ a random action is taken, with probability $1 - \epsilon$ instead the action with greater action-value (Q-value) is taken. Usually, ϵ is progressively decreased during learning.

In Q-Learning the Q-value is computed with respect to a policy (the *target policy*) that would always choose the best action, and never a random one. This is expressed with the formula:

$$Q^*(s, a) = E_{s' \sim p}[R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$$

Equation 4.1

I wrote Q^* instead of Q^π to make explicit that it is not computed over a generic policy π but ideally over the optimal policy (the asterisk means “optimal policy”). You can compare the differences between equations 4.1 and 2.11.

Hence the Q-Value is referred to a tentatively optimal policy, that is asymptotically better or equal than the behavior policy. This makes Q-Learning formally an off-policy method because it evaluates a different policy than the one used to generate trajectories.

In Q-Learning the true Q^* function is unknown, and the algorithm tries to estimate it during the course of learning. So, for the Q-value estimate instead of Q^* we will use the symbol Q_ψ that means that it is an estimate using parameters ψ .

$$Q_\psi(s, a) = E_{s' \sim P} [R(s, a, s') + \gamma \max_{a'} Q_\psi(s', a')]$$

Equation 4.2

The execution of Q-Learning is as follows: at the beginning all $Q_\psi(s, a)$ for all states and actions are set to arbitrary values (if Q_ψ is in tabular form the values can be set close to zero for faster convergence, while if Q_ψ is a function approximator such as a neural network the parameters ψ may be set to a small random noise such as in [Glorot and Bengio 2010] or [Saxe et al. 2013]). Then using the behavior policy (for instance the ϵ -greedy policy specified before) an action a is taken, new state s' is reached, and the reward $r = R(s, a, s')$ is obtained.

Now, if we want our estimate of the Q function to respect the Bellman Equation 4.2, we have to make $Q_\psi(s, a)$ move towards the right-hand side value, to be more similar to it. So, calling the “target value” as y :

$$y = r + \gamma \max_{a'} Q_\psi(s', a')$$

Equation 4.3

Then if Q_ψ is in tabular format, to move it closer to the target value, it is updated in the following way:

$$Q_{\psi'}(s, a) \leftarrow Q_\psi(s, a) + \eta (y - Q_\psi(s, a))$$

Equation 4.4

Where η is the learning rate.

If instead Q_ψ is not in tabular format but it is estimated using a function approximator such as a neural network, a Loss is computed for the discrepancy of the approximated Q-value :

$$L = (y - Q_{\psi}(s, a))^2$$

Equation 4.5

Then the parameters of Q_{ψ} are updated with a gradient descent step (or with batch gradient descent) using the gradient of that Loss function.

It must be noted that in equation 4.5 the value y is to be considered as a constant, not as a value depending on Q_{ψ} , so that the gradient of the loss with respect to ψ will not include the effect of $\max_{a'} Q_{\psi}(s', a')$ but only of $Q_{\psi}(s, a)$. For this reason, the Q-Learning update is considered “semi-gradient”, and not complete gradient descent. In other words, the update in case of function approximation would be the following:

$$\psi \leftarrow \psi + \eta (y - Q_{\psi}(s, a)) \nabla_{\psi} Q_{\psi}(s, a)$$

Equation 4.6

Then, the algorithm continues in the same way: a new action a' is taken following the policy, and Q_{ψ} is updated either with the tabular or the gradient descent method using the new $y' = r' + \gamma \max_{a''} Q_{\psi}(s'', a'')$ as described above, and so on.

Algorithm 1 **Deep Q-Learning Algorithm (simple version)**

Require: Step size η

Require: Initialize parameters ψ of network Q_{ψ} with small random values

Do until (supposed) convergence:

using some behavior policy collect a transition $\langle s_i, a_i, s_{i+1}, r_i \rangle$

$y_i = r + \gamma \max_{a_{i+1}} Q_{\psi}(s_{i+1}, a_{i+1})$

$\psi \leftarrow \psi + \eta (y_i - Q_{\psi}(s_i, a_i)) \nabla_{\psi} Q_{\psi}(s_i, a_i)$

end of do

The algorithm above is a simple version and in the next part of the chapter we will see how to improve it. The first thing to notice is that it collects a transition and computes the respective change of Q-value estimate immediately. But to improve learning speed it is possible to collect a certain number of transitions, called “mini-batch”, and calculate at once all the q-value changes for all the transitions in the mini-batch by vectorized computations in a GPU or in a Tensor Processing Unit (note: in Machine Learning literature some authors use the word

“batch” as a synonym for “mini-batch”, some authors instead reserve the word “batch” to describe the whole set of samples available. Hence for clarity’s sake I used the word “mini-batch”, to indicate a certain number of samples but not the whole set).

As wrote above, Q-Learning is an off-policy algorithm because it trains under a policy that includes exploration, but computes a Q-function estimate for a policy that only assumes exploitation. In fact, because of how the algorithm is designed, theoretically the behavior policy could be even a policy very different from the target policy, we are not obliged to use an ϵ -greedy policy to generate samples. That means that all past trajectories may be retained and reused in future optimization iterations. Nonetheless there is a practical limit: if the q-value estimate is computed through a neural network, it has a finite approximation capacity. A neural network can approximate better when the inputs at inference time are of the same distribution of the training samples (this, very trivially, because the number of parameters is finite and the network cannot retain all the information that is presented during training, and it will approximate better for the kind of inputs that went through a greater number of optimization steps). So, it is better to train the q-network with input states that are one the same distribution that the q-network would experience if using the target policy (the tentatively optimal/greedy policy). To have similar distribution of states it is necessary to use a behavior policy that is similar to the target policy. That motivates the usage of an ϵ -greedy policy and motivates the practice of re-using the old samples only for a certain number of steps and then discard them. Nonetheless, the bigger is your neural network, and hence the network approximation capacity, the longer the old samples may be retained and reused by the algorithm.

Q-learning uses “Bootstrapping” in Q-value updates: it uses Q-value estimate of next state to update the Q-value estimate of current state, instead of using only rewards or complete returns (methods that only use complete returns are instead called “Monte Carlo methods”).

Since Q-Learning is both off-policy and bootstrapping, it is not guaranteed to converge when using function approximators such as neural networks for the Q-function [Szepesvári 2010]. In fact, the three conditions “function approximation”, “bootstrapping”, “off-policy training” are named “The Deadly Triad” when they occur together in value based methods, see [Sutton and Barto 2018, Ch.11]. Despite that, it is believed that Q-Learning may converge if the behavior policy is sufficiently close to the target policy, like in ϵ -greedy policies with small ϵ . In many experiments Q-Learning showed to work well, such as in [Mnih et al. 2013] where a variant of Q-Learning with a convolutional neural network as a function approximator was trained to play Atari video games.

When Q_ψ is a function approximator such as a neural network, the fact that it is used both to compute the prediction (and hence changes at each gradient descent iteration) and the target, makes the bootstrapped value y a “moving target” (literally!). This slows down learning because the gradient descent trajectory is wandering too much. To stabilize the trajectory of the gradient descent it is better to do a mini-batch update with a big number of samples. To further stabilize it, it is useful to observe the fact that it is an off-policy algorithm: each tuple of $\langle s_t, a_t, s_{t+1}, r_t \rangle$ can be saved to a “replay buffer” and used in future: sampling randomly from the replay buffer instead of using the recently experienced trajectory will avoid to having correlated samples that would bias the gradient. In addition to this, to decrease even more the “moving target” effect it would be better to run some cycles of mini-batch updates using a “target network” that is a copy of the q-network but it is not updated at every gradient descent step, and hence it is more stable.

Algorithm 2 Deep Q-Learning Algorithm with Replay Buffer and Target Network

Require: an empty replay buffer B

Require: Step size η

Require: Initialize parameters ψ of network Q_ψ with small random values

do:

copy target network parameters $\psi' \leftarrow \psi$

do N times:

using some behavior policy collect transitions $\{\langle s_i, a_i, s_{i+1}, r_i \rangle\}$, add it to B

do K times:

sample a mini-batch of M transitions $\langle s_i, a_i, s_{i+1}, r_i \rangle$ from B

for each i of M transitions do (all can be done at once if vectorized):

$$y_i = r + \gamma \max_{a_{i+1}} Q_{\psi'}(s_{i+1}, a_{i+1})$$

$$\psi \leftarrow \psi + \eta (y_i - Q_\psi(s_i, a_i)) \nabla_\psi Q_\psi(s_i, a_i)$$

end of for each

end of do

end of do

end of do

Another problem of Q-Learning is that it learns a Q-value that instead of being the maximum of the expectation of the action-value is biased towards the estimate of the maximum of the action-value and hence overestimates the action-value (the maximum of an expectation is different from the expectation of the maximum).

To see how that happens, let us create a toy scenario. Imagine we are in a certain state s_0 in which you have five actions available: a_1, a_2, a_3, a_4, a_5 . Taking a_1 will give a very low reward and lead to a state s_1 with $\max_{a'} Q(s_1, a')$ equal to some constant K . Taking a_2 will give sometimes a very high reward (50% of the times) and sometimes a medium reward (50% of the times) and lead to a state s_2 with a $\max_{a'} Q(s_2, a')$ equal to the same K as s_1 . Actions a_3, a_4 and a_5 will behave in the same way as a_2 , except they will lead respectively to states s_3, s_4 and s_5 , which have the same $\max_{a'} Q(s', a')$ as s_1 and s_2 (that is K). Let us put that in numbers:

Action in s_0	Reward	Leads to State	$\max_{a'} Q(s', a')$
a_1	0 (Probability 1.0)	s_1	$\max_{a'} Q(s_1, a') = K$
a_2	30 (Probability 0.5)	s_2	$\max_{a'} Q(s_2, a') = K$
	10 (Probability 0.5)		
a_3	30 (Probability 0.5)	s_3	$\max_{a'} Q(s_3, a') = K$
	10 (Probability 0.5)		
a_4	30 (Probability 0.5)	s_4	$\max_{a'} Q(s_4, a') = K$
	10 (Probability 0.5)		
a_5	30 (Probability 0.5)	s_5	$\max_{a'} Q(s_5, a') = K$
	10 (Probability 0.5)		

Table 4.1

So we can compute the expected values for each action:

$$\begin{aligned}
 EQ(s_0, a_1) &= 0 * 1.0 + \gamma \max_{a'} Q(s_1, a') = 0 + \gamma K = \gamma K \\
 EQ(s_0, a_2) &= 30 * 0.5 + 10 * 0.5 + \gamma \max_{a'} Q(s_2, a') = 20 + \gamma K \\
 EQ(s_0, a_3) &= 30 * 0.5 + 10 * 0.5 + \gamma \max_{a'} Q(s_3, a') = 20 + \gamma K \\
 EQ(s_0, a_4) &= 30 * 0.5 + 10 * 0.5 + \gamma \max_{a'} Q(s_4, a') = 20 + \gamma K \\
 EQ(s_0, a_5) &= 30 * 0.5 + 10 * 0.5 + \gamma \max_{a'} Q(s_5, a') = 20 + \gamma K
 \end{aligned}$$

Hence, the *maximum of expectation of Q with respect to the actions* for the state s_0 is:

$$\max_{a'} EQ(s_0, a') = 20 + \gamma K$$

But if instead we run the Q-Learning algorithm we will find a different value. For instance imagine we do 15 runs starting from state s_0 , we choose each action three times, and because of random chances we get the following rewards:

Run	Chosen action	Reward r
1	a_1	0
2	a_1	0
3	a_1	0
4	a_2	10
5	a_2	30
6	a_2	30
7	a_3	10
8	a_3	10
9	a_3	10
10	a_4	30
11	a_4	30
12	a_4	30
13	a_5	10
14	a_5	30
15	a_5	10

Table 4.2

Given these outcomes, we can see that randomly, due to the probabilistic nature of rewards, in certain cases the target value $y = r + \gamma \max_{a'} Q_{\psi}(s', a')$ (eq. 4.3) is greater than the expected value $EQ(s_0, a')$ and sometimes it is smaller. For instance for action a_5 the target value is twice $y = 10 + \gamma K$ and only once $y = 30 + \gamma K$. The Q-learning algorithm will then tend to compute a value for $Q(s_0, a_5)$ that is smaller than $EQ(s_0, a_5)$. For action a_4 instead, all three runs obtained the maximum reward, hence the target value is $y = 30 + \gamma K$ all the three times, that is greater than the expected value of $20 + \gamma K$, so the Q-learning algorithm will tend to compute a value for $Q(s_0, a_4)$ that is bigger than $EQ(s_0, a_4)$. At first glance this may seem ok, because we expect that due to random nature of runs, some values will have an estimate greater than their true expectation, and some a lesser one. But the problem is that the Q-Value is the maximum among all the estimates, so if the Q-value estimate for the optimal action (the

one that has actually the greater $EQ(s_t, a')$ is greater than its expected value, it will drive the value of s_t up. But if instead the Q-value estimate for the optimal action is lesser than its expected value, it is not certain that it will drive the value of s_t down, because another action (with expected value lower or equal), due to randomness, may have had an estimate greater than the estimate of the optimal action, and that will still drive the value of s_t up. The value of a state (computed as the maximum among its Q-values) is used to bootstrap the Q-values of the states that reach it, propagating the value overestimation to other states.

In our example we know that :

$$\max_a EQ(s_0, a) = EQ(s_0, a_2) = EQ(s_0, a_3) = EQ(s_0, a_4) = EQ(s_0, a_5)$$

We also know that

$$Q_\psi(s_0, a_2) > EQ(s_0, a_2)$$

$$Q_\psi(s_0, a_3) < EQ(s_0, a_3)$$

$$Q_\psi(s_0, a_4) > EQ(s_0, a_4)$$

$$Q_\psi(s_0, a_5) < EQ(s_0, a_5)$$

$$Q_\psi(s_0, a_4) > Q_\psi(s_0, a_2)$$

And we know that $V(s_0)$ is estimated by:

$$\max_a Q_\psi(s_0, a) = Q_\psi(s_0, a_4) > \max_a EQ(s_0, a)$$

Hence even if for optimal actions we had an equal number of runs with target greater and with target lesser than expected, the estimate of s_0 value $\max_a Q_\psi(s_0, a)$ is greater than the real value, and through bootstrapping will propagate the overestimation to any other Q-values of states that reach s_0 .

So in this toy example we see how random fluctuations may give returns that are either lesser or greater than true average returns but while all times the return of optimal action is greater than its expectation it drives the $\max_a Q_\psi(s, a)$ up, not all times the return of optimal action is lesser than its expectation it drives the $\max_a Q_\psi(s, a)$ down. This creates an overestimation bias that propagates to other Q-values through bootstrapping.

A formal mathematical proof of the overestimation, as well as the proposal of an algorithm named “*Double Q-Learning*” that mitigates the problem can be found in [van Hasselt 2010]. In Double Q-Learning the algorithm uses two different action-value function estimators, e.g. two neural networks that we may call A and B. When computing the target value, one neural

network (let us say A) is used to check which action has the maximum value, while the other network (let us say B) is used to actually obtain the Q-value of that action. Then the target value is used to update the Q-value of the former network (A, in this case).

At each iteration the estimators may swap roles: the one that previously was used to evaluate which action has the maximum value may now be used to obtain the Q-value of the optimal action, while the estimator who was previously used to get the Q-value of the optimal action may now be used to find which action has the maximum value, and be subject to Q-value update. The roles of the two networks A and B may be decided randomly or in a fixed way after a certain number of iterations. For clarity's sake in the pseudocode algorithm we use the names "network Q_{ψ_1} " and "network Q_{ψ_2} " instead of A and B.

Algorithm 3 Deep Double Q-Learning Algorithm with Replay Buffer and Target Network

Require: an empty replay buffer B

Require: Step size η

Require: Initialize parameters ψ_1 of network Q_{ψ_1} with small random values

Require: Initialize parameters ψ_2 of network Q_{ψ_2} with small random values

do:

 copy target network parameters $\psi_1' \leftarrow \psi_1$

 copy target network parameters $\psi_2' \leftarrow \psi_2$

 do N times:

 using some behavior policy collect transitions $\{ \langle s_i, a_i, s_{i+1}, r_i \rangle \}$, add it to B

 do K times:

 sample a mini-batch of M transitions $\langle s_i, a_i, s_{i+1}, r_i \rangle$ from B

 choose (e.g. randomly) either UPDATE(1) or UPDATE(2)

 for each i of M transitions do (all can be done at once if vectorized):

 if UPDATE(1) then:

$$a_{i+1}^* = \underset{a_{i+1}}{\arg \max} Q_{\psi_1'}(s_{i+1}, a_{i+1})$$

$$y_i = r + \gamma Q_{\psi_2'}(s_{i+1}, a_{i+1}^*)$$

$$\psi_1 \leftarrow \psi_1 + \eta (y_i - Q_{\psi_1}(s_i, a_i)) \nabla_{\psi_1} Q_{\psi_1}(s_i, a_i)$$

 else if UPDATE(2) then:

$$a_{i+1}^* = \underset{a_{i+1}}{\arg \max} Q_{\psi_2'}(s_{i+1}, a_{i+1})$$

$$y_i = r + \gamma Q_{\psi_1'}(s_{i+1}, a_{i+1}^*)$$

$$\psi_2 \leftarrow \psi_2 + \eta (y_i - Q_{\psi_2}(s_i, a_i)) \nabla_{\psi_2} Q_{\psi_2}(s_i, a_i)$$

 end if

 end of for each

 end of do

 end of do

end of do

It has to be noted that the Q-Learning algorithms presented above can be used only with a set of discrete actions. The problem with continuous actions is the computation of the maximum q-value among all possible actions: when we have a set of discrete actions it is possible to check the q-value of all actions and pick the maximum, but if the actions are continuous we

would need an analytic way to compute the maximum of the q-function, and that is not possible if the q-function estimate is represented by a neural network. There have been many attempts to use workarounds to extend Q-Learning to continuous actions, but for a matter of space we will not describe them here. A better way to do Reinforcement Learning with continuous actions is to use Policy Based methods, that are detailed in the following chapter.

5. Policy Based Methods

Policy based methods, also called “Policy Gradient methods” and “Policy Optimization methods” (those two expressions have diverse meaning but often are used interchangeably), differently from the value based methods, may or may not use value or action-value functions, and are characterized instead by having a parametrized policy function that, given the state/observation as input, returns a probability for each action, without explicitly assigning expected returns values.

In other terms, the policy function $\pi(\cdot | s_t)$ is a function that is explicitly parametrized by θ (for instance θ may be the weights of a neural network), such that if at time t the state is s_t , it will return the probability distribution over which action to take a_t :

$$a_t \sim \pi_\theta(\cdot | s_t)$$

Equation 5.1

That is like equation 2.1 with the addition of parameters θ .

One way in which a neural network can output a probability distribution over a set of discrete actions is to have a last layer with as many neurons as the number of actions, and on top of that operate a Softmax. If instead the actions are continuous then the last layer of the neural network will have as many neurons as continuous parameters in the action, and each neuron will be considered the mean of the distribution of that action parameter (the distribution type may be arbitrary, for instance a gaussian usually works well).

The problem of Reinforcement Learning is still the same, we want to maximize the expected return $J(\pi)$, (detailed in equation 2.5), but this time we would like to modify directly the policy

parameters θ . One way of doing that could be to do gradient ascent using the gradient of the expected return $J(\pi)$ with respect to the parameters θ :

$$\theta_{k+1} \leftarrow \theta_k + \eta \nabla_{\theta} J(\pi_{\theta_k})$$

Equation 5.2

With η learning rate.

So it is necessary to compute the gradient of the expected return $J(\pi_{\theta_k})$ with respect to θ .

Apparently that could seem problematic because the expected return $J(\pi_{\theta_k})$ is an expectation with respect to the distribution of the trajectories, that depends both on the policy and on the transition function, but we assume to not know the transition function so we do not know the distribution of trajectories (or the distribution of states). Hence, we do not even know how a change of the policy may change the distribution of trajectories. But we need to estimate how the expected return changes as a consequence of changes of policy (due to changes of θ), and this seems to be problematic because, as just said, the effect of the policy change on the trajectories distribution is unknown. Fortunately, as we will see in next chapter, there is a way to derive the gradient of the expected return $J(\pi_{\theta_k})$ that does not involve the effect of the change of policy on states/trajectories distribution, so this is not really a problem.

6. Policy Gradient

A basic algorithm that improves the policy through gradient ascent on the expected return $J(\pi_{\theta})$ with respect to the policy parameters θ is the one called “**Vanilla Policy Gradient**” [OpenAi 2018A] (or just “Policy Gradient”), that I describe below. It is very similar to the earlier “REINFORCE” [Williams 1992], with the difference that REINFORCE updates the policy parameters at each step while Vanilla Policy Gradient updates parameters using mini-batches of steps. It is an on-policy algorithm. Some use the term “vanilla policy gradient” as a family of algorithms, with REINFORCE being one of them (see [Peters and Schaal 2008]), where the policy gradient is not manipulated (as it happens instead in Natural Policy Gradient or in Proximal Policy Optimization, which I describe in later chapters).

Following the method in [OpenAi 2018A] it is useful to recall from calculus that the derivative of $\log(x)$ with respect to x is $1/x$. Then, because of chain rule, the derivative of $\log(f(x))$

with respect to x is the derivative of $f(x)$ divided by $f(x)$. Applying that to the trajectory distribution, and rearranging, we obtain the so-called “*log derivative trick*”:

$$\nabla_{\theta} P(\tau|\theta) = P(\tau|\theta) \nabla_{\theta} \log P(\tau|\theta)$$

Equation 6.1

Now, taking the definition of the probability of a trajectory τ made of T steps from equation 2.4, we compute the log of it:

$$\log P(\tau|\theta) = \log \rho_0(s_0) + \sum_{t=0}^{T-1} \log P(s_{t+1}|s_t, a_t) + \log \pi_{\theta}(a_t|s_t)$$

Equation 6.2

At this point we want to compute $\nabla_{\theta} \log P(\tau|\pi)$. We note that the terms $P(s_{t+1}|s_t, a_t)$ and $\rho_0(s_0)$ are environmental and that do not depend on θ , so their gradient is zero. Hence:

$$\nabla_{\theta} \log P(\tau|\theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$$

Equation 6.3

Now, computing the gradient of equation 2.5 we obtain:

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} E_{\tau \sim \pi_{\theta}} [G(\tau)] \\ &= \nabla_{\theta} \int_{\tau} P(\tau|\theta) G(\tau) \\ &= \int_{\tau} \nabla_{\theta} P(\tau|\theta) G(\tau) \end{aligned}$$

Equation 6.4

And applying equation 6.1 (log derivative trick):

$$= \int_{\tau} P(\tau|\theta) \nabla_{\theta} \log P(\tau|\theta) G(\tau)$$

$$= E_{\tau \sim \pi_\theta} [\nabla_\theta \log P(\tau|\theta) G(\tau)]$$

Equation 6.5

Applying equation 6.3 :

$$\therefore \nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) G(\tau) \right]$$

Equation 6.6

The same formula can be obtained also in a different way through the *Policy Gradient Theorem* [Sutton et al. 1999], [Sutton and Barto 2018, Ch.13] which I will not detail here.

We obtained the gradient of the expected return in form of an expectation. This means that if we can sample what is contained inside the expectation, we can use it to estimate the gradient. It turns out that we can do that: the expectation is with respect to the trajectory τ , whose probability distribution is determined by the environment and by our policy π_θ , so by making the agent follow the policy we can collect the returns $G(\tau)$ and compute the estimate of the expectation. Since the agent must follow the policy π_θ the algorithm is on-policy .

If we run N different episodes or trajectories $\{\tau_0, \tau_1, \dots, \tau_{N-1}\}$, with respective trajectory lengths $\{T_0, T_1, \dots, T_{N-1}\}$ we can do mini-batch gradient ascent. If we call the actions and states of trajectory i at time t respectively with $a_{i,t}$ and $s_{i,t}$ the following is mean gradient:

$$\hat{g} = \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) G(\tau_i)$$

Equation 6.7

Then we can update the weights of the policy neural network:

$$\theta' \leftarrow \theta + \eta \hat{g}$$

Equation 6.8

Algorithm 4 Policy Gradient

Require: Policy network step size η

Require: Initialize parameters θ of network π_θ with small random values

For $k = 0, 1, 2, \dots$ do:

 collect trajectories $D_k = \{ \tau_i \}$ using policy $\pi_k(\theta_k)$

 compute rewards-to-go $G(\tau_t)$

 estimate policy gradient as:

$$\widehat{g}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) G(\tau_t)$$

 update the policy network with gradient ascent (or other methods like Adam):

$$\theta_{k+1} \leftarrow \theta_k + \eta \widehat{g}_k$$

end of for

A notable fact is that the gradient of the log probability of each action $a_{i,t}$ is multiplied by the complete return $G(\tau_i)$. The complete return of the trajectory includes the rewards that have been received before that action, and for which the action has not any responsibility. From a logical standpoint we would expect that the gradient of the log probability of a certain action is multiplied only by the sum of the rewards obtained after that action, instead also the rewards obtained before that action are used. As we will see, the part of gradient of the expected return that involves the multiplication for the rewards obtained before the action is equal to zero in expectation, so we can use only the rewards obtained after the action. To prove it, I must digress a little and introduce the Expected Grad-Log-Prob (EGLP) lemma, from [OpenAI 2018A].

EGLP Lemma: suppose that f_θ is a parametrized distribution over a random variable x , then:

$$E_{x \sim f_\theta} [\nabla_\theta \log f_\theta(x)] = 0$$

Equation 6.9

Proof (noting that the integral of a probability distribution is always = 1 by definition):

$$\nabla_\theta \int_x f_\theta(x) = \nabla_\theta 1 = 0$$

Equation 6.10

Now, applying the “log derivative trick” equation 6.1 :

$$\begin{aligned}
0 &= \nabla_{\theta} \int_x f_{\theta}(x) \\
&= \int_x \nabla_{\theta} f_{\theta}(x) \\
&= \int_x f_{\theta}(x) \nabla_{\theta} \log f_{\theta}(x) \\
\therefore 0 &= E_{x \sim f_{\theta}} [\nabla_{\theta} \log f_{\theta}(x)]
\end{aligned}$$

Equation 6.11

(Optional: this proof is analogous to the proof that the expected value of statistical score of a distribution conditioned on the true value of its parameters is equal to zero. In fact since the statistical score is the gradient of the log-likelihood, when θ are the true values of the parameters of the distribution we could substitute the likelihood function $\mathcal{L}(\theta; X)$ in place of $f_{\theta}(x)$ in equation 6.9 , that means that the probability of X would be $\mathcal{L}(\theta; X)$, and we will obtain the same result, that in this case would mean that expected value of statistical score conditioned on true θ is zero).

We can decompose $G(\tau)$ in 2 sums, one of the rewards before the action and one after the action happened at time t (this proof is elaborated from the one by [Soemers 2019], there exists also another proof by [OpenAi 2018B]).

We use $G(\tau_{0:t-1})$ for the past rewards:

$$G(\tau_{0:t-1}) = \sum_{k=0}^{t-1} \gamma^k r_k$$

Equation 6.12

And we call $G(\tau_t)$ “reward-to-go”.

$$G(\tau_t) = \sum_{k=t}^{T-1} \gamma^k r_k$$

Equation 6.13

$$G(\tau) = G(\tau_{0:t-1}) + G(\tau_t)$$

Equation 6.14

We can then decompose equation 6.6 in the same way:

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_{0:t-1}) + \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_t) \right] \\ &= E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_{0:t-1}) \right] + E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_t) \right] \end{aligned}$$

Equation 6.15

Now, this can be rewritten as:

$$part1 = E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_{0:t-1}) \right]$$

Equation 6.16

$$part2 = E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G(\tau_t) \right]$$

Equation 6.17

$$\nabla_{\theta} J(\pi_{\theta}) = part1 + part2$$

Equation 6.18

Now, because of EGLP lemma:

$$E_{a_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] = 0$$

Equation 6.19

Now, from 6.16 we have:

$$\begin{aligned}
part1 &= E_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) G(\tau_{0:t-1}) \right] \\
&= \sum_{t=0}^{T-1} E_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) G(\tau_{0:t-1})]
\end{aligned}$$

Equation 6.20

We know that $\nabla_\theta \log \pi_\theta(a_t | s_t)$ and $G(\tau_{0:t-1})$ are independent because the action depends only on the state at time t by definition, and $G(\tau_{0:t-1})$ depends only on states and actions before time t . So, we can separate the expectation of the multiplication into a multiplication of expectations:

$$= \sum_{t=0}^{T-1} E_{\tau \sim \pi_\theta} [G(\tau_{0:t-1})] E_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t)]$$

Equation 6.21

Now we plug equation 6.19 in 6.21:

$$\begin{aligned}
&= \sum_{t=0}^{T-1} E_{\tau \sim \pi_\theta} [G(\tau_{0:t-1})] * 0 \\
&\therefore part1 = 0
\end{aligned}$$

Equation 6.22

So we obtain:

$$\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) G(\tau_t) \right]$$

Equation 6.23

$$\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \sum_{l=t}^{T-1} \gamma^l r_l \right]$$

Equation 6.24

This permits to use, for each action, only the rewards obtained after it (the rewards-to-go) to calculate the gradient of the expected return. Using the complete returns would not be wrong, because it is equal in expectation, but it would introduce much more variance in the gradient and that would slow down learning. Using $\gamma = 1$ would assure that each action has the same importance, it can be used in finite-horizon (episodic learning). In infinite time horizon it is not possible to collect the whole return (the trajectory never ends), so in that case it is necessary to bootstrap $G(\tau_t)$ using only the reward for the first step (or for some steps) and then summing a discounted estimate of the state value.

For instance, in infinite-horizon, a k -step bootstrap, using a function approximator $V_\psi(s)$, parametrized by ψ , to estimate the value function, would be:

$$\widehat{G}_k(\tau_t) = \sum_{l=t}^{t+k-1} \gamma^l r_l + \gamma^{t+k} V_\psi(s_{t+k})$$

Equation 6.25

$\widehat{G}_k(\tau_t)$ is equal in expectation to $G(\tau_t)$, so it is theoretically correct to use it, but since in any algorithm we are not using the real state value but an approximation, it will introduce bias (but it will allow to work with infinite-horizon cases).

In infinite horizon case moreover, since we don't use the whole infinite trajectory at once but only n -step of it a time, we have to consider the infinite trajectory as equivalent to an infinite set of n -step episodes. So, every n -step we will reset t to 0, that is also necessary to avoid that the discount makes progressively less relevant the subsequent rewards. Or, if we don't want to reset t to 0, just change the equation a little, subtracting t to the exponent:

$$\widehat{\widehat{G}}_k(\tau_t) = \sum_{l=t}^{t+k-1} \gamma^{l-t} r_l + \gamma^k V_\psi(s_{t+k})$$

Equation 6.26

In this way the finite and infinite horizon are included in the same framework.

Also $Q^{\pi_\theta}(s_t, a_t)$ is theoretically correct to be used instead of $G(\tau_t)$, because since inside the expectation the actions are distributed following the policy π_θ , it implies that $Q^{\pi_\theta}(s_t, a_t)$ is equal in expectation to $G(\tau_t)$ (another formal proof by OpenAI may be found in [OpenAI 2018C]).

In fact, it is theoretically correct also to use, instead of $G(\tau_t)$, any other expression that is equal in expectation. Every expression that starts with $G(\tau_t)$ or $\widehat{G}_n(\tau_t)$ and then adds any other expression that does not depend on the current action (but may depend on current state, since the current action and any function depending only on the current state are conditionally independent given the current state), is equal in expectation. That is because of EGLP lemma, and we already used it to show that $G(\tau_{0:t-1})$ may be or may be not added to $G(\tau_t)$ without changing the expectation (see how we obtained equation 6.23).

So, more generally:

$$\nabla_{\theta} J(\pi_{\theta}) = E_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \Phi_t \right]$$

Equation 6.27

With Φ_t that may be one of:

$$\begin{aligned} &G(\tau_t) \\ &\widehat{G}_k(\tau_t) \\ &\widehat{\widehat{G}}_k(\tau_t) \\ &Q^{\pi_{\theta}}(s_t, a_t) \\ &G(\tau_t) + b(s_t) \\ &\widehat{G}_k(\tau_t) + b(s_t) \\ &\widehat{\widehat{G}}_k(\tau_t) + b(s_t) \\ &Q^{\pi_{\theta}}(s_t, a_t) + b(s_t) \end{aligned}$$

Equation 6.28

Where $b(s_t)$, usually called “baseline”, is an additional expression that may depend on state and may be negative. For example, it may be $V^{\pi_{\theta}}(s_t)$ (in practice we would use its estimate $V_{\psi}(s_t)$).

So, many different Φ_t are possible. For instance if we start from $Q^{\pi_{\theta}}(s_t, a_t)$ and we add $b(s_t) = -V^{\pi_{\theta}}(s_t)$, we end up with the Advantage function $\Phi_t = A^{\pi_{\theta}}(s_t, a_t) = Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)$, see equation 2.13 .

The advantage function makes sense because the policy gradient update is intended to increase the probability of actions that perform better than the policy and to decrease the

probability of actions that perform worse than the policy, if we use the Advantage function as Φ_t we are sure to have positive updates when the action is better than current policy choice, and negative updates when it is worse. This will have the effect of reducing variance in the gradient updates and hence speed up learning.

In case of using an advantage function, it is necessary to have a neural network $V_\psi(s)$, parametrized by ψ , that estimates the value function. That neural network will be updated with a mini-batch gradient descent, similar to the Q-network gradient descent of equations 4.3 and 4.4, where the target value y may be computed by $G(\tau_t)$, the Monte Carlo rewards-to-go of state t (this would have a great variance) or may be computed bootstrapping the value using $V_\psi(s')$ (in that case it will not be full gradient descent but semi-gradient, since the value of the next state will be considered as a constant and not as a function of ψ).

The value network update will be the following:

$$\begin{aligned}
 y &= r + \gamma V_\psi(s') \text{ if bootstrapping} \\
 &\text{or} \\
 y &= G(\tau_t) = \sum_{k=t}^{T-1} \gamma^k r_k \text{ if using Monte Carlo} \\
 L &= (y - V_\psi(s))^2 \\
 \psi &\leftarrow \psi + \eta (y - V_\psi(s)) \nabla_\psi V_\psi(s)
 \end{aligned}$$

Equation 6.29

Policy gradient methods that use an estimate of the value function to bootstrap the estimate of the total return such in equation 6.25 and 6.26 are called “**actor-critic**”, where actor is referred to the policy function and critic to the value function [Sutton and Barto 2018, Ch.13]. Some authors call “actor-critic” every policy method that uses an estimate of the value function as component of Φ_t , but Sutton and Barto specify that the term should be exclusive for those that use value function to bootstrap the estimate of the total return. That means when generally, for a fixed k and a certain baseline $b(s_t)$:

$$\Phi_t = \sum_{l=t}^{t+k-1} \gamma^{l-t} r_l + \gamma^k V_\psi(s_{t+k}) - b(s_t)$$

Equation 6.30

It is common to do that with $k = 1$ and with the estimate of the value function used as baseline, so that:

$$\Phi_t = r_t + \gamma V_\psi(s_{t+1}) - V_\psi(s_t) ,$$

Equation 6.31

Which is equivalent to a bootstrapped advantage function.

Algorithm 5 Policy Gradient with rewards-to-go and advantage function

Require: Policy network step size η

Require: Value network step size ω

Require: Initialize parameters θ of network π_θ with small random values

Require: Initialize parameters ψ of network V_ψ with small random values

For $k = 0, 1, 2, \dots$ do:

 collect trajectories $D_k = \{ \tau_i \}$ using policy $\pi_k(\theta_k)$

 compute rewards-to-go $G(\tau_t)$

 compute advantage estimates \widehat{A}_t using current estimate of value function V_{ψ_k} :

$$\widehat{A}_t = G(\tau_t) - V_{\psi_k}(s_t)$$

 estimate policy gradient as:

$$\widehat{g}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \widehat{A}_t$$

 update the policy network with gradient ascent (or other methods like Adam):

$$\theta_{k+1} \leftarrow \theta_k + \eta \widehat{g}_k$$

For $n = 0, 1, 2, \dots, N-1$ do a value function gradient descent iteration:

 estimate the value function gradients as:

$$\widehat{h}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \frac{1}{T} \sum_{t=0}^{T-1} \nabla_\psi (V_\psi(s_t) - G(\tau_t))^2$$

 update the value function with a gradient descent step (or other method):

$$\psi_{k+1} \leftarrow \psi_k + \omega \widehat{h}_k$$

 end of for

end of for

It must be noted that in this algorithm the policy ascent direction is not optimal, we will describe the problem and present an algorithm that computes a better direction in Ch. 8 with Natural Policy Gradient.

One thing to be careful of in Policy Gradient is the fact that when the policy step size η is too little, the learning will be slow. So it should be big enough to make the learning progress quick, but not so big to make the new policy worse than the old. In fact the new policy may end up being quite close in parameters space to the old policy, but even small differences in parameters may result in big differences in policies. Hence a small step in parameter space may correspond to a big step in policy space, in a way that worsens the policy instead of improving it. An algorithm that aims at computing the right step size is Trust Region Policy Optimization, in ch. 9.

As it is evident watching the Policy Gradient algorithm, each trajectory is used only for one step of optimization in the policy network gradient ascent. This is because it is an on-policy algorithm: since the gradient of the expected returns (equation 6.27) is an expectation with respect to the distribution of samples generated by the policy π_θ the generated samples can be used only once, to improve π_θ . They cannot be used again to improve the next policy $\pi_{\theta'}$ (the new policy resulted from the optimization step) because they come from a different distribution, the one generated by the old policy. It may be tempting to use each trajectory more than once, but as [Schulman et al. 2017] report: *"doing so is not well-justified, and empirically it often leads to destructively large policy updates"*. A method that instead allows to use each trajectory more than once is Proximal Policy Optimization, described in Ch. 10, which permits to use the samples generated by the previous policy distribution if the new optimized policy distribution is not too different.

Another way to re-use the samples generated by a previous policy, or more generally a way to use samples generated by a different policy (that means off-policy learning) is to use importance sampling as described by [Degris et al. 2012], [Levine and Koltun 2013] and [Levine 2021b]. The problem with importance sampling is that some data points that would occur rarely with the behavior policy may produce very big gradients, because of the multiplication of the ratio of the two policy probabilities that in this case would be high, and this may cause instability. This problem is not existent in the Deterministic Policy Gradient algorithm [Silver et al. 2014], which is an off-policy policy gradient method that does not use importance sampling.

A remarkable aspect of policy gradient methods is that compared to value based methods they have a more "natural" distribution on actions: the policy function outputs "by design" a distribution on actions, and during training the distribution smoothly changes towards one that

improves the returns. In value based methods instead the estimate of value function or q-function only gives expected (action-)state returns, so we have to choose an artificial way to create a distribution on actions, such as with ϵ -greedy. Hence in value base methods when some action that previously seemed not optimal starts to seem better than all the other actions (that is when its q-value estimate becomes the greater among all actions), we have an abrupt change in the policy: the policy stops suggesting the previously-considered-best action (except for a small probability depending on ϵ) and starts suggesting almost greedily the new best action. In policy gradient methods instead the change in distribution is smooth because it is due to gradual changes through the policy gradient ascent, and this also implies an automatic gradual passage from exploration to exploitation.

7. Generalized Advantage Estimation

[Schulman et al. 2016] introduced a particular version of the Advantage Function, the “Generalized Advantage Estimation”, with proven variance reduction property. The setting is the policy gradient family of algorithms, using undiscounted rewards. So there is not a time-dependent parameter γ to discount the rewards, but there is another parameter named γ with a different meaning but same mathematical behaviour: it is meant to downweight rewards corresponding to delayed effects. Thus, the meaning is different from the “discount”, it does not mean that delayed effects are less “useful” at current time (that would be the meaning of the classical time-discount), but it still discounts the delayed effects, and doing so it introduces bias, because in this way the total return will be smaller in absolute value. On the other hand, it will decrease variance, for the same reason: smaller absolute value means smaller variance among different trajectories. So, the meaning of this parameter $\gamma \in [0,1]$ here is to be a switch for smoothly parametrize a trade-off between bias and variance, with bias increasing and variance decreasing when γ decreases, and bias decreasing and variance increasing when γ tends to 1.

Then this Generalized Advantage Estimation has another characteristic that decreases variance but introduces bias: it computes the estimate of the returns as an exponentially weighted average of k different n -step bootstrap estimates, where n goes from 1 to k . To be clearer, if we define:

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

Equation 7.1

Imagine of having a trajectory of length at least $k>1$, then we could build a 1-step bootstrap estimate of the advantage function at time t :

$$\hat{A}_t^{(1)} = r_t + \gamma V(s_{t+1}) - V(s_t) = \delta_t^V$$

Equation 7.2

But, if $k>2$ we could also build a 2-step bootstrap estimate, or if $k>3$ a 3-step estimate, or a k -step estimate etc. :

$$\hat{A}_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) - V(s_t) = \delta_t^V + \gamma \delta_{t+1}^V$$

$$\hat{A}_t^{(3)} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 V(s_{t+3}) - V(s_t) = \delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V$$

...

$$\hat{A}_t^{(k)} = \sum_{l=t}^{t+k-1} \gamma^{l-t} r_l + \gamma^k V(s_{t+k}) - V(s_t) = \sum_{l=t}^{t+k-1} \gamma^{l-t} \delta_l^V$$

$$\hat{A}_t^{(k)} = \widehat{\widehat{G}}_k(\tau_t) - V(s_t)$$

Equation 7.3

In infinite horizon, k may tend to ∞ . In that case the $\gamma^k V(s_{t+k})$ at each time will be canceled by the discounted $-V(s_t)$ of the next time and we would have:

$$\hat{A}_t^{(\infty)} = \sum_{l=t}^{\infty} \gamma^{l-t} r_l - V(s_t) = \sum_{l=t}^{\infty} \gamma^{l-t} \delta_l^V$$

Equation 7.4

Now, a 1-step estimate has low variance but high bias, while a 5-step estimate has bigger variance and lower bias, the more steps we include in the estimate the bigger the variance and the lower the bias. It is possible to compute an exponentially weighted average of all n -steps estimators, parametrized by λ such that when $\lambda = 0$ the weighted average coincides with $\hat{A}_t^{(1)}$ and when $\lambda = 1$ it coincides with $\hat{A}_t^{(\infty)}$, so that λ is another parameter that may be used to tune the bias/variance trade-off.

To obtain that, let us note that the series $\sum_{k=0}^{\infty} \lambda^k$, with $\lambda \in [0,1]$, is equal to $1/(1-\lambda)$. Hence $(1-\lambda) \sum_{k=0}^{\infty} \lambda^k = 1$. So, a summation of infinite terms, each weighted by $(1-\lambda)\lambda^k$, has the total sum of weights = 1.

The “Generalized Advantage Function” hence is:

$$\hat{A}_t^{GAE(\gamma, \lambda)} = (1 - \lambda) \sum_{k=1}^{\infty} \lambda^{k-1} \hat{A}_t^{(k)}$$

Equation 7.5

And a more practical formulation may be obtained:

$$\begin{aligned} \hat{A}_t^{GAE(\gamma, \lambda)} &= (1 - \lambda)(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots) \\ &= (1 - \lambda)(\delta_t^V + \lambda(\delta_t^V + \gamma \delta_{t+1}^V) + \lambda^2(\delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V) + \dots) \\ &= (1 - \lambda)(\delta_t^V(1 + \lambda + \lambda^2 + \lambda^3 + \dots) + \gamma \delta_{t+1}^V(\lambda + \lambda^2 + \lambda^3 + \lambda^4 \dots) \\ &\quad + \gamma^2 \delta_{t+2}^V(\lambda^2 + \lambda^3 + \lambda^4 + \lambda^5 \dots) + \dots) \\ &= (1 - \lambda)(\delta_t^V \left(\frac{1}{1 - \lambda} \right) + \gamma \delta_{t+1}^V \left(\frac{\lambda}{1 - \lambda} \right) + \gamma^2 \delta_{t+2}^V \left(\frac{\lambda^2}{1 - \lambda} \right) + \dots) \\ &= \sum_{l=t}^{\infty} (\lambda \gamma)^{l-t} \delta_l^V \end{aligned}$$

Equation 7.6

This GAE hence can be used instead of the Advantage Function and put in place of Φ_t in the policy gradient (see equation 6.27). This permits to tune variance and bias with two parameters: γ controls how far in time to consider the rewards relevant, and λ that decides how much the return component of the advantage has to be similar to a 1-step bootstrap (if $\lambda=0$) or progressively to a Monte Carlo return (if $\lambda=1$).

8. Natural Policy Gradient

To improve the speed of training one right thing to do can be to scale the gradient in a way that gradient descent works better. It is possible that the gradient is not pointing to the direction of steepest ascent, and it may be convenient to modify or scale the gradient differently for each parameter in order to make it point to the direction of greater policy improvement

[Kakade 2002], [Peters et al. 2003], similarly to what has been found in supervised learning [Amari 1998].

In fact, some parameters of the neural network affect the policy more than others: when we are doing gradient ascent with a certain learning rate, we basically are putting a limit on how much parameters can be changed, but this is not equivalent to put a limit on how much the policy can be changed. A faster learning would happen if we were able to fix the rate at which the policy is changed, and have larger rates for the parameters with little influence on the policy and smaller rates for the parameters with greater influence on the policy.

One way to do that consists of putting a constraint to each gradient descent step, such that the new policy and the old policy are not too different as distributions. An appropriate computation for that could be the Kullback-Leibler divergence (see Appendix A.1), enforcing the constraint $D_{KL}(\pi_{\theta'}, || \pi_{\theta}) < \epsilon$.

The Kullback-Leibler divergence is defined as $D_{KL}(\pi_{\theta'}, || \pi_{\theta}) = E_{\theta'}[\log \pi_{\theta'} - \log \pi_{\theta}]$.

Second order Taylor expansion of KL divergence is approximated by (see Appendix A.4 and A.5):

$$D_{KL}(\pi_{\theta'}, || \pi_{\theta}) \approx \frac{1}{2}(\theta' - \theta)^T \mathbf{F}_{\pi_{\theta}}(\theta' - \theta)$$

Equation 8.1

Where $\mathbf{F}_{\pi_{\theta}}$ is the Fisher-information matrix of the policy π_{θ} .

The Fisher-information matrix of a distribution is the expected value of the covariance matrix of the score, given parameters θ . The score is the gradient of log-likelihood with respect to the parameters (see Appendix A.2 and A.3). Hence:

$$\mathbf{F}_{\pi_{\theta}} = E_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^T]$$

Equation 8.2

That can be approximated with samples from the trajectory:

$$\widehat{\mathbf{F}}_{\pi_{\theta}} = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(a_i|s_i) \nabla_{\theta} \log \pi_{\theta}(a_i|s_i)^T$$

Equation 8.3

The common gradient ascent step for policy gradient is:

$$\theta' \leftarrow \theta + \eta \nabla_{\theta} J(\pi_{\theta})$$

Equation 8.4

That can be seen as optimizing a constrained first order Taylor expansion of $J(\pi_{\theta})$:

$$\theta' \leftarrow \arg \max_{\theta'} \nabla_{\theta} J(\pi_{\theta})^T (\theta' - \theta)$$

$$\text{such that } \|\theta' - \theta\|^2 \leq \epsilon$$

Equation 8.5

That would imply:

$$\theta' \leftarrow \theta + \sqrt{\frac{\epsilon}{\|\nabla_{\theta} J(\pi_{\theta})\|^2}} \nabla_{\theta} J(\pi_{\theta})$$

Equation 8.6

But as I wrote before, instead of imposing a constraint over $\|\theta' - \theta\|^2$ it would be better to impose a constraint over the distributions, such as $D_{KL}(\pi_{\theta'}, \pi_{\theta}) < \epsilon$.

Given the relation between $D_{KL}(\pi_{\theta'}, \pi_{\theta})$ and Fisher-information matrix from equation 8.1, and given that we can compute an approximation of $\mathbf{F}_{\pi_{\theta}}$ from samples (as in eq. 8.3), such a constraint in the KL divergence can be imposed in a gradient descent step with the following update [Peters and Schaal 2008], [Levine 2021]:

$$\theta' \leftarrow \theta + \eta \mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta})$$

$$\text{with } \eta = \sqrt{\frac{2\epsilon}{\nabla_{\theta} J(\pi_{\theta})^T \mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta})}}$$

Equation 8.7

This is called “Natural gradient” (or “Covariant Gradient” as in [Bagnell and Schneider 2003]), and its usage permits a faster learning with a smoother gradient ascent.

The computation of η as in eq. 8.7 is necessary only if you want to adhere to a strict a-priori fixed ϵ , otherwise you can use an arbitrarily small value for η : doing so will not change the direction of the change.

The derivation of η as in eq. 8.7 can be found in Appendix A.1 of [Peters 2007] and we propose it here in an equivalent form in the following lines.

Let us define $\delta = (\theta' - \theta)$. So we can redefine the optimization problem, starting from eq. 8.5, using the Kullback-Leibler constraint, and the eq. 8.1 approximation, as:

$$\delta \leftarrow \arg \max_{\delta} \nabla_{\theta} J(\pi_{\theta})^T \delta$$

$$\text{such that } D_{KL}(\pi_{\theta+\delta} || \pi_{\theta}) < \epsilon$$

Equation 8.8

This is a constrained optimization problem, solvable by Lagrangian Multiplier method. The resulting Lagrangian is the following:

$$\Lambda(\delta, \lambda) = \nabla_{\theta} J(\pi_{\theta})^T \delta + \lambda (\epsilon - \frac{1}{2} \delta^T \mathbf{F}_{\pi_{\theta}} \delta)$$

Equation 8.9

Now, all partial derivatives of $\Lambda(\delta, \lambda)$ should be zero.

$$\nabla_{\delta} \Lambda(\delta, \lambda) = \nabla_{\theta} J(\pi_{\theta}) - \lambda \frac{1}{2} 2 \mathbf{F}_{\pi_{\theta}} \delta = 0$$

$$\nabla_{\theta} J(\pi_{\theta}) = \lambda \mathbf{F}_{\pi_{\theta}} \delta$$

$$\lambda^{-1} \mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta}) = \lambda^{-1} \mathbf{F}_{\pi_{\theta}}^{-1} \lambda \mathbf{F}_{\pi_{\theta}} \delta$$

$$\delta = \lambda^{-1} \mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta})$$

Equation 8.10

Now we plug the δ found in eq. 8.10 into the constraint $\frac{1}{2} \delta^T \mathbf{F}_{\pi_{\theta}} \delta = \epsilon$ and we have the dual function:

$$\frac{1}{2} \lambda^{-1} (\mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta}))^T \mathbf{F}_{\pi_{\theta}} \lambda^{-1} \mathbf{F}_{\pi_{\theta}}^{-1} \nabla_{\theta} J(\pi_{\theta}) = \epsilon$$

$$\frac{1}{2\epsilon} (\mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta))^T \nabla_\theta J(\pi_\theta) = \lambda^2$$

because of matrix properties $AB^T = B^T A^T$

$$\lambda^2 = \frac{1}{2\epsilon} \nabla_\theta J(\pi_\theta)^T (\mathbf{F}_{\pi_\theta}^{-1})^T \nabla_\theta J(\pi_\theta)$$

since \mathbf{F}_{π_θ} is symmetric, also $\mathbf{F}_{\pi_\theta}^{-1}$ is symmetric, and for symmetric matrices $A = A^T$

$$\lambda = \sqrt{\frac{\nabla_\theta J(\pi_\theta)^T \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)}{2\epsilon}}$$

Equation 8.11

Wit λ the Lagrange multiplier. Now we plug it into eq. 8.10:

$$\delta = \sqrt{\frac{2\epsilon}{\nabla_\theta J(\pi_\theta)^T \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)}} \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)$$

Equation 8.12

Since δ is the name we used for $(\theta' - \theta)$, we know that:

$$\theta' - \theta = \sqrt{\frac{2\epsilon}{\nabla_\theta J(\pi_\theta)^T \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)}} \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)$$

$$\theta' = \theta + \sqrt{\frac{2\epsilon}{\nabla_\theta J(\pi_\theta)^T \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)}} \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)$$

$$\text{now if we name } \eta = \sqrt{\frac{2\epsilon}{\nabla_\theta J(\pi_\theta)^T \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)}} \text{ we have:}$$

$$\theta' = \theta + \eta \mathbf{F}_{\pi_\theta}^{-1} \nabla_\theta J(\pi_\theta)$$

Equation 8.13

The final 2 lines of eq. 8.13 are equivalent to eq. 8.7 .

9. Trust Region Policy Optimization (incomplete)

In the Trust Region Policy Optimization an alternative optimization problem is posed, such that allows to maximize the length of the step of the policy gradient ascent.

The Reinforcement Learning problem here is framed in a different formulation.

In the classical Policy Gradient formulation, we aim at maximizing $J(\pi_\theta)$, that for ease of reading we may now write $J(\theta)$. But equivalently, we may maximize the performance of a policy (the one to be maximized) with respect to a fixed policy (the one used until now to generate trajectories). If we call $\pi_{\theta'}$ the optimized policy (parametrized by the new parameters θ') and π_θ the fixed policy, we may want to maximize $J(\theta') - J(\theta)$.

Following [Schulman et al. 2015], it is useful to note that:

$$J(\theta') - J(\theta) = E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t A^{\pi_\theta}(s_t, a_t) \right]$$

Equation 9.1

Proof:

Recall equation 2.12 and 2.13:

$$Q^{\pi_\theta}(s_t, a) = E_{s_{t+1} \sim P} [R(s_t, a, s_{t+1}) + \gamma V^{\pi_\theta}(s_{t+1})]$$

$$A^{\pi_\theta}(s_t, a) = Q^{\pi_\theta}(s_t, a) - V^{\pi_\theta}(s_t)$$

Then:

$$\begin{aligned} & E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t A^{\pi_\theta}(s_t, a_t) \right] \\ &= E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t (R(s_t, a, s_{t+1}) + \gamma V^{\pi_\theta}(s_{t+1}) - V^{\pi_\theta}(s_t)) \right] \\ &= E_{\tau \sim \pi_{\theta'}} \left[-V^{\pi_\theta}(s_0) + \sum_t \gamma^t R(s_t, a, s_{t+1}) \right] \end{aligned}$$

$$= E_{\tau \sim \pi_{\theta'}}[-V^{\pi_{\theta}}(s_0)] + E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t R(s_t, a, s_{t+1}) \right]$$

distribution of s_0 does not depend on policy but only on initial state distribution $\rho_0 \Rightarrow$

$$= E_{s_0 \sim \rho_0}[-V^{\pi_{\theta}}(s_0)] + E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t R(s_t, a, s_{t+1}) \right]$$

$$= E_{\tau \sim \pi_{\theta}}[-V^{\pi_{\theta}}(s_0)] + E_{\tau \sim \pi_{\theta'}} \left[\sum_t \gamma^t R(s_t, a, s_{t+1}) \right]$$

$$= -J(\theta) + J(\theta') \quad \therefore$$

Equation 9.2

Now, equation 9.1 means that the difference between the expected return of the policy parametrized by θ' and the expected return of the policy parametrized by θ is equal to the expectation of the advantage function $A^{\pi_{\theta}}(s_t, a_t)$ over trajectories distributed by the policy parametrized by θ' . (Since the expectation of the trajectories follows the distribution by policy θ' , it implies that the actions a_t are distributed by $\pi_{\theta'}$, so $A^{\pi_{\theta}}(s_t, a_t)$ in eq. 9.1 it is the advantage function of the actions selected by the new policy θ' with respect to the fixed policy θ).

We can improve the policy if we can maximize the right-hand side of eq. 9.1 (because the greater is the difference between the expected returns of new policy and the expected returns of the old policy, the greater is the right-hand side). To do that, one strategy is to have an equivalent equation in a form that we can sample, so that we can run gradient ascent with respect to the policy parameters θ .

We can define the probability of being in state s at time t , depending on policy parametrized by θ as $Prob(s_t = s \mid \theta)$.

Then we can define the frequency $\xi_{\theta}(s)$ of a state s as the number of times that s is expected to be visited, computed as unnormalized (it is a frequency, not a probability) and time-discounted, under the policy π_{θ} parametrized by θ :

$$\xi_{\theta}(s) = Prob(s_0 = s | \theta) + \gamma Prob(s_1 = s | \theta) + \gamma^2 Prob(s_2 = s | \theta) + \dots$$

Equation 9.3

We can rewrite eq. 9.1 in a way to sum over time, then over states, and then over actions. Then in a second passage we can write it in a way to sum over states first, and in the last passage we plug eq. 9.3, referred to policy θ' , into it:

$$J(\theta') - J(\theta) = \sum_t \sum_s Prob(s_t = s | \theta') \sum_a \pi_{\theta'}(a_t | s_t) \gamma^t A^{\pi_{\theta}}(s_t, a_t)$$

$$J(\theta') - J(\theta) = \sum_s \sum_t \gamma^t Prob(s_t = s | \theta') \sum_a \pi_{\theta'}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t)$$

$$J(\theta') - J(\theta) = \sum_s \xi_{\theta'}(s) \sum_a \pi_{\theta'}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t)$$

Equation 9.4

This equation entails that any new policy parametrized by θ' which in every state has a better or equal expected advantage function with respect to the previous policy parametrized by θ , i.e. $\sum_a \pi_{\theta'}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t) \geq 0$, improves the total expected return (or leaves the total expected return unchanged if the expected advantage function is zero at each state).

Now, the right-hand side seems in a more manageable form to be maximized, but if we look carefully, we see that to sample it we would need to follow the frequency of states $\xi_{\theta'}(s)$, referred to the new policy θ' , but at that point we only know the old policy θ and we can sample only with that. Hence, we use a local approximation that uses the old policy for the frequencies (please note the usage of $\xi_{\theta}(s)$ instead of $\xi_{\theta'}(s)$):

$$L_{\xi}(\theta, \theta') = \sum_s \xi_{\theta}(s) \sum_a \pi_{\theta'}(a_t | s_t) A^{\pi_{\theta}}(s_t, a_t)$$

Equation 9.5

This approximation can be used only when the new distribution $\pi_{\theta'}(a_t | s_t)$ is not too different from the old distribution $\pi_{\theta}(a_t | s_t)$, it must stay inside a “trust region” (hence the name of the algorithm). We will see later how that is mathematically translated.

We need also to express this formula as an expectation to sample it.

Any series $\sum_{k=0}^{\infty} \gamma^k$ with $\gamma \in [0,1]$ is equal to $1/(1 - \gamma)$.

Hence, the sum over the frequencies $\sum_s \xi_{\theta}(s)$ can then be replaced by the surrogate expectation $\frac{1}{1-\gamma} E_{s \sim \xi_{\theta}(s)}$.

Since the formula will be used for optimization, we can get rid of the constant $\frac{1}{1-\gamma}$.

The sum over actions can be replaced by an expectation over actions, and we want those actions to be generated by the old policy θ , while now they are multiplied by the probabilities of the new policy θ' , so we need to apply importance sampling. So, the new objective to be maximized becomes:

$$L_{\xi}(\theta, \theta') = E_{s \sim \xi_{\theta}(s)} \left[E_{a \sim \pi_{\theta}(a|s)} \left[\frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)} A^{\pi_{\theta}}(s, a) \right] \right]$$

Equation 9.6

Now, if we look at equation 9.6 we note that the expectation with respect to the states $E_{s \sim \xi_{\theta}(s)}$ is with respect to the discounted frequency of the states $\xi_{\theta}(s)$ (e.q. 9.3), that even when normalized through a multiplication by $1/(1 - \gamma)$ is not generally exactly the same as the state distribution that happens during training, which is determined by the transition function $P(s'|s, a)$ and the policy π_{θ} , or equivalently by the trajectory distribution $P(\tau|\pi_{\theta})$ (eq. 2.4). This is because of the presence of γ^t term in $\xi_{\theta}(s)$ (unless in case of some non-general assumptions, such as when the visiting probability of each state never changes, being the same at each time t).

But when we sample our trajectories we actually sample from $P(\tau|\pi_{\theta})$.

So, we need to write a slightly different surrogate objective, using the expectation that is consistent with our sampling, obtaining the following function:

$$\begin{aligned} L(\theta, \theta') &= E_{s \sim P(\tau|\pi_{\theta})} \left[E_{a \sim \pi_{\theta}(a|s)} \left[\frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)} A^{\pi_{\theta}}(s, a) \right] \right] \\ &= E_{\tau \sim \pi_{\theta}} \left[\frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)} A^{\pi_{\theta}}(s, a) \right] \end{aligned}$$

Equation 9.7

Fortunately, it turns out that we can use that as objective. This is justified by the fact that the gradient of equation 9.7 matches the gradient of $J(\theta')$ if evaluated locally at the point $\theta' = \theta$:

$$\begin{aligned}
\nabla_{\theta'} L(\theta, \theta') \Big|_{\theta' = \theta} &= E_{\tau \sim \pi_{\theta}} \left[\frac{\nabla_{\theta'} \pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)} A^{\pi_{\theta}}(s, a) \right] \Big|_{\theta' = \theta} \\
&= E_{\tau \sim \pi_{\theta}} [\log \nabla_{\theta'} \pi_{\theta'}(a|s) A^{\pi_{\theta}}(s, a)] \Big|_{\theta' = \theta} \\
&= E_{\tau \sim \pi_{\theta}} [\log \nabla_{\theta'} \pi_{\theta'}(a|s) A^{\pi_{\theta}}(s, a)] \Big|_{\theta' = \theta} \\
&= \nabla_{\theta'} J(\theta') \Big|_{\theta' = \theta}
\end{aligned}$$

Equation 9.8

So this means that a small change in the policy parameters from θ to θ' such that $L(\theta, \theta')$ increases, makes also the new policy's expected return $J(\theta')$ increase with respect to the old policy expected return $J(\theta)$, but nothing is said about what exactly is this “small change”.

Now we need to use the concept of Kullback-Leibler divergence so the reader who does not know about it may read Appendix A.1 before proceeding.

Considering the status s where the Kullback-Leibler divergence occurring between $\pi_{\theta}(\cdot | s)$ and $\pi_{\theta'}(\cdot | s)$ is maximum, we define the value of such divergence as D_{KL}^{max} :

$$D_{KL}^{max}(\pi_{\theta}, \pi_{\theta'}) = \max_s [D_{KL}(\pi_{\theta}(\cdot | s) || \pi_{\theta'}(\cdot | s))]$$

Equation 9.9

Then, we define the maximum absolute value that the Advantage function can possibly assume, in any possible pairs of state and action, as ϵ :

$$\epsilon = \max_{s,a} |A^{\pi_{\theta}}(s, a)|$$

Equation 9.10

Then it can be proved, following formula 9 in TRPO paper [Schulman et al. 2017], that

$$J(\theta') - J(\theta) \geq L(\theta, \theta') - C D_{KL}^{max}(\pi_{\theta}, \pi_{\theta'})$$

$$\text{where } C = \frac{4 \epsilon \gamma}{(1 - \gamma)^2}$$

Equation 9.11

We may consider C as a penalty coefficient.

Now, if we call θ_i the policy parameters at time i , and we call ω the policy parameters for any other different policy, we define:

$$M_i(\omega) = J(\theta_i) + L(\theta_i, \omega) - C D_{KL}^{max}(\pi_{\theta_i}, \pi_{\omega})$$

Equation 9.12

Then, if we start from eq. 9.11, substituting θ' with ω and θ with θ_i , and rearranging a little, we have:

$$\begin{aligned} J(\omega) &\geq J(\theta_i) + L(\theta_i, \omega) - C D_{KL}^{max}(\pi_{\theta_i}, \pi_{\omega}) \\ J(\omega) &\geq M_i(\omega) \end{aligned}$$

Equation 9.13

Now, if $\omega = \theta_i$ it is easy to see that since the Kullbach-Leibler divergence of two identical distribution is zero, and the advantage function of two identical policies is zero (that makes $L(\theta_i, \theta_i) = 0$):

$$J(\theta_i) = M_i(\theta_i)$$

Equation 9.14

Now, let us say that at time $i + 1$ we have a policy parametrized by θ_{i+1} . This can be the case such as in a loop in which at each iteration we change the policy parameters. By eq. 9.14 we have:

$$J(\theta_{i+1}) \geq M_i(\theta_{i+1})$$

Equation 9.15

We can combine 9.14 with 9.15:

$$J(\theta_{i+1}) - J(\theta_i) \geq M_i(\theta_{i+1}) - M_i(\theta_i)$$

Equation 9.16

This implies that maximizing M_i at each iteration implies maximizing the expected returns !

Now, if you watch $M_i(\omega)$ in the eq. 9.12, you notice that $J(\theta_i)$ is fixed, so you only need to maximize $L(\theta_i, \omega) - C D_{KL}^{max}(\pi_{\theta_i}, \pi_\omega)$, that means maximizing L including a penalty that depends on how much the new policy distribution differs from the old policy distribution.

As [Schulman et al. 2017] notice, using that C penalty coefficient is theoretically justified but leads to very small steps.

Alternatively, we could be able to take larger steps if we just maximize L and put as a condition that the difference between the new policy distribution and the old policy distribution stays under a certain threshold δ (the “trust region”), measuring that difference with a Kullback-Leibler divergence.

$$\underset{\theta'}{\text{maximize}} L(\theta, \theta') = E_{\tau \sim \pi_\theta} \left[\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)} A^{\pi_\theta}(s, a) \right]$$

$$s. t. \quad D_{KL}^{max}(\pi_\theta, \pi_{\theta'}) \leq \delta$$

Equation 9.17

This constraint on the KL divergence applies to all points in space and it is not practical to be concretely applied, so an heuristic approximation may be used instead, computing the average Kullback-Leibler difference (which at computation time will be the average of the samples).

$$\underset{\theta'}{\text{maximize}} L(\theta, \theta') = E_{\tau \sim \pi_\theta} \left[\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)} A^{\pi_\theta}(s, a) \right]$$

$$s. t. \quad E_{\tau \sim \pi_\theta} [D_{KL}(\pi_\theta(\cdot | s) || \pi_{\theta'}(\cdot | s))] \leq \delta$$

Equation 9.18

At this point the optimization problem is completely defined. We need only an effective way to solve it. The original authors [Schulman et al. 2015] used $Q^{\pi_\theta}(s, a)$ instead of $A^{\pi_\theta}(s, a)$, which leads to an equivalent optimization problem because it changes the objective only by a constant, but saves from the burden of computing $V^{\pi_\theta}(s)$.

$$\underset{\theta'}{\text{maximize}} \quad L(\theta, \theta') = E_{\tau \sim \pi_\theta} \left[\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)} Q^{\pi_\theta}(s, a) \right]$$

$$\text{s.t.} \quad E_{\tau \sim \pi_\theta} [D_{KL}(\pi_\theta(\cdot|s) || \pi_{\theta'}(\cdot|s))] \leq \delta$$

Equation 9.19

$Q^{\pi_\theta}(s, a)$ is computed by an empirical estimate, such as the sum of discounted rewards of the sampled trajectory that starts with (s, a) , that is actually $G(\tau)$ of eq. 2.3 .

The parameters θ are then updated with the θ' found solving the constrained optimization problem, with the conjugate gradient algorithm followed by line search (point 3 of Ch. 6, and Appendix C of [Schulman et al. 2015]).

The objective $L(\theta, \theta')$ is approximated by its first order Taylor expansion $\hat{L}(\theta, \theta')$. Since in $L(\theta, \theta')$ the parameters θ are fixed and the variable parameters are θ' , and the Taylor expansion is built around the point $\theta' = \theta$, we have that:

$$\hat{L}(\theta, \theta') = L(\theta, \theta) + \nabla_{\theta'} L(\theta, \theta')^T \Big|_{\theta' = \theta} (\theta' - \theta)$$

$$\text{since } L(\theta, \theta) = 0 \Rightarrow$$

$$\hat{L}(\theta, \theta') = L(\theta, \theta')^T \Big|_{\theta' = \theta} (\theta' - \theta)$$

we plug eq. 9.8 in

$$\hat{L}(\theta, \theta') = \nabla_{\theta'} J(\theta')^T \Big|_{\theta' = \theta} (\theta' - \theta)$$

Equation 9.20

As explained in Appendix A.6, the Hessian of Kullback-Leibler divergence between two distributions of the same family $\pi_\theta(x)$ and $\pi_{\theta'}(x)$, with respect to θ' , evaluated at $\theta' = \theta$ is equal to the Fisher Information Matrix of $\pi_\theta(x)$. Hence this can be used to approximate the Fisher Information Matrix: using the Hessian of Kullback-Leibler divergence instead of standard $E_{x \sim P_\theta(x)} [\nabla_\theta \log P_\theta(x) \nabla_\theta \log P_\theta(x)^T]$ (as in Appendix A.3).

$$\widehat{\mathbf{F}}_{\mathbf{P}_\theta}[i, j] = \frac{1}{N} \sum \frac{\partial^2}{\partial \theta'_i \partial \theta'_j} D_{KL}(\pi_\theta(\cdot | s) || \pi_{\theta'}(\cdot | s))$$

Equation 9.21

To compute the Hessian of Kullback-Leibler divergence of $\pi_\theta(a, s)$ and $\pi_{\theta'}(x)$, for each sampled state s we need to take in consideration all possible actions a , not only the one that has been actually taken by the sampled trajectory. An analytic estimator must integrate over a for each sampled state s .

The Kullback-Leibler divergence is computed by the approximation described in Appendix A.4 (eq. a.18).

$$\widehat{D}_{KL}(f(x; \theta) || f(x; \theta + \delta)) = \frac{1}{2} \delta^T \widehat{\mathbf{F}}_{\mathbf{P}_\theta} \delta$$

Equation 9.22

The original authors [Schulman et al. 2015] experimented also a sampling scheme named "Vine", applicable on any environment in which it is possible to restart from a certain well-defined state (as in a simulated environment). The Vine sampling consists of trying different actions from the same starting point, generating different trajectories to compute an average Q estimate, which in this way will have a lower variance.

<< to be completed in next version >>

10. Proximal Policy Optimization

The Proximal Policy Optimization algorithm [Schulman et al. 2017] builds on the ideas and theoretical framework of Trust Region Policy Optimization. To understand this chapter it is necessary to have read the TRPO chapter (Ch. 9) until equation 9.18 included. In fact, it aims at optimizing the same eq. 9.18:

$$\underset{\theta'}{\text{maximize}} \quad L(\theta, \theta') = E_{\tau \sim \pi_\theta} \left[\frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)} A^{\pi_\theta}(s, a) \right]$$

$$\text{s.t.} \quad E_{\tau \sim \pi_\theta} [D_{KL}(\pi_\theta(\cdot | s) || \pi_{\theta'}(\cdot | s))] \leq \delta$$

While in TRPO there was only one policy gradient ascent iteration in each optimization round (in which the step length was computed so to not make the new policy too different from the old), in PPO the algorithm does more than one policy gradient ascent iteration in each optimization round using the same set of trajectories: the trajectories are re-used as long as the new policy is not too different from the old.

What distinguishes Proximal Policy Optimization from TRPO is that it uses a different strategy to constraint the KL divergence: (1) it uses gradient clipping as a way to never have a too big gradient update, and (2) it checks when the old and new policies are too divergent so to stop reusing the samples generated by the old policy.

For ease of reading, let us call the ratio between the two probabilities $d(\theta) = \frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)}$. A new surrogate objective function that uses clipping, called L^{CLIP} may be devised:

$$L^{CLIP}(\theta, \theta') = E_{\tau \sim \pi_{\theta}} [\min (d(\theta)A^{\pi_{\theta}}(s, a) , \text{clip}(d(\theta), 1 - \epsilon, 1 + \epsilon) A^{\pi_{\theta}}(s, a))]$$

Equation 10.1

written in a simplified way:

$$L^{CLIP}(\theta, \theta') = E_{\tau \sim \theta} [\min (d(\theta)A^{\pi_{\theta}} , \text{clip}(d(\theta), 1 - \epsilon, 1 + \epsilon) A^{\pi_{\theta}})]$$

Equation 10.2

That means that, generating actions with policy θ , the objective is the expectation of the minimum among $d(\theta)A^{\pi_{\theta}}$ and $\text{clip}(d(\theta), 1 - \epsilon, 1 + \epsilon) A^{\pi_{\theta}}$.

The *clip* function imposes both a ceiling and a floor to $d(\theta)$: if $d(\theta) > 1 + \epsilon$ it returns $1 + \epsilon$, if $d(\theta) < 1 - \epsilon$ it returns $1 - \epsilon$, otherwise it returns $d(\theta)$.

The min function then has the result of returning a lower bound on the unclipped objective. In this way, when the probability ratio would make the objective improve, it is bounded to be $< 1 + \epsilon$, so to not have too big steps. While when it would make the objective worse it is unbounded towards negative infinity. When the objective is negative, it is bounded towards zero so to have always some negative value that has an impact on the policy update.

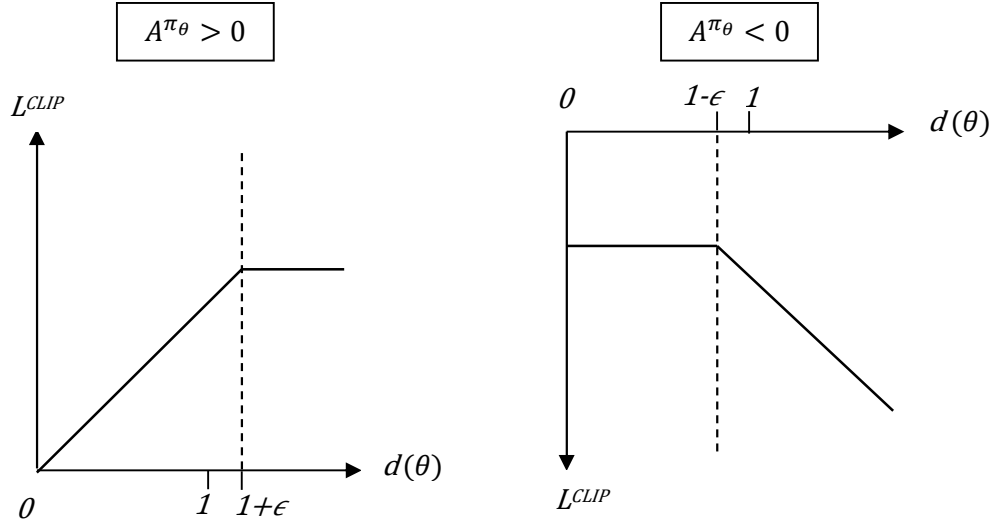


Figure 10.1 Surrogate objective with policy ratio clipping in Proximal Policy Optimization

To practically use it, we have to do gradient descent with respect to θ' on equation 10.1 or 10.2 , using an automatic differentiation library that supports clipping.

This clipping does not guarantee that the new policy is not too different from the old policy, so it is necessary to check the Kullback-Leibler divergence between the probabilities of the taken actions under the old policy and under the new policy, and stop the policy gradient ascent iterations in case it is over a certain threshold (“early stopping”). This avoids using samples generated from a policy that has a distribution too different from the one that gets optimized. The KL divergence computed in this way is just an approximation because is the average of the KL divergence for each sample. Since it is sampled on the distribution of actions from the old policy, it is already the empirical expectation over the old policy (so there is no need to multiply for $\pi_{\theta_{old}}(a_t|s_t)$ in the divergence formula):

$$\widehat{Kl} = \frac{1}{T} \sum_t \log (\pi_{\theta_{old}}(a_t|s_t)) - \log (\pi_{\theta}(a_t|s_t))$$

Equation 10.3

Algorithm 6 Proximal Policy Optimization with ratio clipping and early stopping

Require: Policy network step size η

Require: Value network step size ω

Require: Threshold for Kullback-Leibler divergence based early stopping ζ

Require: Initialize parameters θ of network π_{θ} with small random values

Require: Copy parameters θ_{old} of network $\pi_{\theta_{old}}$ from θ

Require: Initialize parameters ψ of network V_{ψ} with small random values

For $k = 0, 1, 2, \dots$ do:

collect trajectories $D_k = \{ \tau_i \}$ using policy $\pi_{\theta_{old}}$

compute rewards-to-go $G(\tau_t)$

compute advantage estimates \widehat{A}_t using current estimate of value function V_{ψ_k} :

$$\widehat{A}_t = G(\tau_t) - V_{\psi_k}(s_t)$$

For $m = 0, 1, 2, \dots, M - 1$ do a policy gradient ascent iteration:

compute policy ratios:

$$d_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$$

compute clipped surrogate losses:

$$L_t^{CLIP}(\theta, \theta') = \min(d_t(\theta)\widehat{A}_t, \quad clip(d_t(\theta), 1 - \epsilon, 1 + \epsilon) \widehat{A}_t)$$

compute KL divergence between $\pi_{\theta_{old}}(a_t|s_t)$ and $\pi_{\theta}(a_t|s_t)$ on taken actions:

$$\widehat{KL} = \frac{1}{|D_k|} \sum_{\tau \in D_k} \frac{1}{T} \sum_{t=0}^{T-1} \log(\pi_{\theta_{old}}(a_t|s_t)) - \log(\pi_{\theta}(a_t|s_t))$$

If $\widehat{KL} \geq \zeta$:

stop doing policy gradient ascent and exit this internal For-cycle

end of If

estimate policy gradient as:

$$\widehat{g}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \sum_{t=0}^{T-1} \nabla_{\theta} L_t^{CLIP}(\theta, \theta')$$

update the policy network with gradient ascent (or other methods like Adam):

$$\theta \leftarrow \theta + \eta \widehat{g}_k$$

end of for

Copy optimized policy into fixed policy: $\theta_{old} \leftarrow \theta$

For $n = 0, 1, 2, \dots, N - 1$ do a value function gradient descent iteration:

estimate the value function gradients as:

$$\widehat{h}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \frac{1}{T} \sum_{t=0}^{T-1} \nabla_{\psi} (V_{\psi}(s_t) - G(\tau_t))^2$$

update the value function with a gradient descent step (or other method):

$$\psi_{k+1} \leftarrow \psi_k + \omega \widehat{h}_k$$

end of for

end of for

PPO is considered on-policy because it aims at optimizing the same policy function that generates the training trajectories (and it is not valid to use an arbitrary policy function to generate trajectories). But actually, it is also *slightly* off-policy because it does possibly more than one step of optimization with the same trajectories, and at each step further than the first it is optimizing a policy that is not *exactly* the same that generated trajectory: it is very similar (small KL-divergence) but not the same.

11. Reward Shaping and Curriculum Learning

Reinforcement Learning may be defined as a method to automatically learn from experience without human help. In fact, this is not completely correct. One of the crucial aspects of actual Reinforcement Learning algorithms is the modelling of the rewards: it is up to the human designer to decide when a reward should be given and how high (or low) the reward is, and that must be done in a way that reflects the actual goal that the agent is supposed to learn. Reward modelling can be made in a simple way if the problem has a single, unique goal: a positive reward is given if the goal is reached, no other rewards are given while the goal is not reached. While this seems to make sense, it may work badly: the agent may reach a state very close to goal without fulfilling the goal, so it would not receive any reward for it and it would not learn that the reached state was good and desirable to be reached again. Since RL algorithms “propagate back” the information obtained from a state that obtained reward to the previous states that lead to it, agents are not able to learn when a state is ideally good if there is not a reward deriving from that state or future reached states. So, it is clear that for some tasks, having only a final reward when the goal is reached is not the best way to assign rewards, and the reward system should instead be designed and engineered. This means that the human choices about how to model the rewards may change a lot the performance of the same algorithm.

Designing the reward system is called “Reward Shaping” and it has to be done in a way that “guides” the agent into intermediate subgoals. It must be noted that Reward Shaping is somehow like “cheating”: the more you shape rewards, the more you are inserting human knowledge in a system that was meant to learn automatically, without prior knowledge. That does not mean it is a bad thing, it just means that we have to be aware that we are inserting external knowledge, and potentially any bias or flaw that comes with it.

In more complicated settings, where the goal may not be unique and there are different “good things to do”, there may be rewards for minor goals, and sometimes this may lead to the fact that the agent learns only to solve the minor goals because the obtained rewards distract it from solve the bigger goal (small rewards may be much smaller than big rewards but much easier to get).

Another method used to make agents learn in complicated task is the “curriculum learning”: a simplified version of the task (or the environment) is used at the beginning, and when the agent has learnt to reach the goal in that setting, a progressively more complete version is used.

Curriculum learning may also be intended as making the agent learn a set of certain skills or behaviours that are preparatory for the complete goal.

Also curriculum learning, like reward shaping, is injecting human knowledge into the system, because a human is deciding which simplified task or skill is necessary to be learnt before learning the true task, and that opens the possibility to bias the learning (i.e. without a certain curriculum learning maybe the agent would learn a better policy that is not based on what the curriculum designer thought were the necessary skills).

12. Imitation Learning

At this point it may be due a brief digression about Imitation Learning, a method that, in common with Reinforcement Learning, aims at selecting the best action depending on the state. Differently from RL, Imitation Learning does not use reinforcement signals such as rewards, but rather it is a supervised learning method in which examples reproduce the behaviour of an expert (that usually is a human). The input data X represent the state, and the label Y is the action taken by the expert in that state. Imitation Learning has the valuable characteristic of being able to use directly human knowledge in form of examples. Unfortunately, one of the issues of IL is that training examples are usually created in a limited subset of the state space, and a trained agent may, during his functioning, exit from that subspace, because even small differences in action responses or in starting states may accumulate and lead to very different states. Once the agent is outside the subspace of states in which training examples have been produced (outside the training set distribution), it likely will not be able to generalize the new states, and the actions selected will not be optimal, or

may even be disastrous. For this reason, it is important to cover a big part of the state space with training samples, and that could be difficult.

An Imitation Learning system may be built as a neural network, with a number of input neurons equal to the dimension of a sample, a certain number of hidden layers with a variety of possible architectures, then a last layer whose output neurons are the ones that indicate the action to take (so if the actions are discrete and there are K different actions there will be K final neurons with a Softmax applied to them, if the actions are continuous there will be other K output neurons outputting the average of the computed distributions on the continuous values, etc.). In other words, it will be a neural network just like the one that you would build for a policy network in RL, except that this network will be trained as a supervised network with usual (mini-batch) gradient descent algorithm, and not with any policy gradient algorithm.

A first thing to note is that the fact that the policy network may have the same architecture both in imitation learning and in policy gradient RL makes possible to train the same network both with IL and RL, so it is theoretically possible to integrate human knowledge with the RL process.

An interesting characteristic of Imitation Learning training is its mathematical relationship with Reinforcement Learning, in particular with policy gradient methods, as we will see it below.

Let's start with the example of discrete actions. To use the same notation of Reinforcement Learning, we call $\pi_\theta(a_t|s_t)$ the Imitation Learning neural network that outputs a probability distribution over actions taking the state as input, and we imagine of having N example trajectories of length T_i each.

Now, the last layer of the network is a Softmax which implies that the gradient of the Loss is the negative of the log of the probability of the correct action. So, it turns out that this Imitation Learning mini-batch gradient (of Loss function) will be computed as:

$$\nabla_\theta J_{IL}(\pi_\theta) = -\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t})$$

Equation 11.1

While in case of Reinforcement Learning the mini-batch gradient (of expected returns) is (from eq. 6.7 and eq. 6.27):

$$\nabla_{\theta} J_{RL}(\pi_{\theta}) = \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \Phi_{i,t}$$

Equation 11.2

Remarkably, they are similar (except for the initial minus in IL because in IL we minimize the objective while in RL we maximize the objective), the Reinforcement Learning gradient is just like the Imitation Learning gradient in which each sample gradient is also multiplied by $\Phi_{i,t}$, that is related to the rewards (it may be the rewards-to-go, the Advantage function etc.).

The same happens if the actions are not discrete but continuous, represented by a real number for each action parameter, so the output of the network is not a Softmax operator but real numbers, intended to be the means of a distribution, usually a Gaussian. For simplicity let us deal with the case of just one action parameter. In that configuration the Loss function of the Imitation Learning system would be the mean square error between the network output and the parameter chosen by the expert (from the sample). The RL formula for the log probability of the policy network will again be similar to the one of IL. To see that, for the RL policy network let us call σ^2 the variance of the gaussian of the action parameter and call $\pi_{\theta}(s_t)$ the output of the policy network, whose value is to be intended as the mean of the gaussian distribution of the parameter of the action. It is necessary to carefully not confuse the policy output $\pi_{\theta}(s_t)$ with π , the transcendental number π that is necessary for the gaussian formula and appearing in the equation. The agent will select the value of action a_t by sampling from a gaussian distribution with mean $\pi_{\theta}(s_t)$ and variance σ^2 . So, the probability of the action taken by the agent will be the probability of the action a_t considering that it is distributed by that gaussian.

$$\begin{aligned} \nabla_{\theta} J_{RL}(\pi_{\theta}) &= \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\pi_{\theta}(s_{i,t}) - a_{i,t})^2 / (2\pi\sigma^2)} \Phi_{i,t} = \\ &= \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} \left(-\frac{(\pi_{\theta}(s_{i,t}) - a_{i,t})^2}{2\pi\sigma^2} - \log \sqrt{2\pi\sigma^2} \right) \Phi_{i,t} = \\ &= -\frac{1}{2\pi\sigma^2} \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} (\pi_{\theta}(s_{i,t}) - a_{i,t})^2 \Phi_{i,t} \end{aligned}$$

Equation 11.3

Now, $1/(2\pi\sigma^2)$ can be cancelled from the equation because it is just a constant that may be incorporated with the learning rate. The same would apply to the Imitation Learning gradient if we derived it within the maximum log-likelihood framework: a $1/(2\pi\sigma^2)$ constant would appear, and we would incorporate it with the learning rate. So, the gradient for Reinforcement Learning policy is:

$$\nabla_{\theta} J_{RL}(\pi_{\theta}) = -\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} (\pi_{\theta}(s_{i,t}) - a_{i,t})^2 \Phi_{i,t}$$

Equation 11.4

While the gradient for the Imitation Learning is given by the usual supervised regression loss gradient, the gradient of the square error:

$$\nabla_{\theta} J_{IL}(\pi_{\theta}) = \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^{T_i-1} \nabla_{\theta} (\pi_{\theta}(s_{i,t}) - a_{i,t})^2$$

Equation 11.5

Again, the two formulas are almost the same, apart from the minus sign (depending on ascending or descending the gradient) and the multiplication by $\Phi_{i,t}$.

So, intuitively, Reinforcement Learning policy gradient is like an Imitation Learning gradient that also uses the information of the rewards $\Phi_{i,t}$ to inform “how good” or “how bad” the action is.

I think it is surprising to find out that two different optimization problems related to two different underlying tasks end up with such similar computations.

13. Afterword

This has been just an introduction to the theory of Deep Reinforcement Learning. To not make it too boring, some mathematical passages and proofs have been skipped, some others have been simplified, but still some have been necessarily reported in detail, to create a consistent presentation that follows a principled thread. This introduction has focused on “Deep” RL, that means that I have not dealt with tabular methods, even if there has been a great amount of

research on them and they cannot be ignored. Even among Deep RL methods, since this is just an introduction, I did not examine all algorithms, but only a didactically representative small subset of them.

The reader who wants the complete math and proofs, as well as the reader that wants to know more about tabular methods or about other deep RL methods, is left to the reference literature.

Also, while algorithms are described in detail, I acknowledge that it may not be easy for the primer to understand exactly how to implement a RL system after reading this introduction. To that purpose I suggest the reader to follow one of the many introductory courses on the internet (there are excellent ones both from MOOCs platforms and from famous brick and mortar universities), that focus on the practical side. A valid help to understand RL algorithms may come from checking the open source implementations of the algorithms available online, some of which are made by the same authors of the reference literature.

A topic that I did not discuss but it is worth a final mention is the fact the Reinforcement Learning is not “sample efficient”, that means that it needs a lot of samples (a lot of experience, or trajectories, or “trial and error”) to learn good policies in complex environments. Often training time may be much longer than expected, for instance longer than (usually) with supervised learning systems. This and other difficulties of Reinforcement Learning practice are well detailed in an article written by [Irpan 2018], which I suggest every RL practitioner to read.

Appendix A: Information Theory Refresh

A.1 Kullback-Leibler Divergence

It is a way to compute how different two distributions are. The Kullback-Leibler divergence between two distributions P and Q (also named *relative entropy* of P from Q) is defined as:

$$\begin{aligned} \text{continuous} \quad D_{KL}(P||Q) &= \int_x P(x) \log \frac{P(x)}{Q(x)} \\ \text{discrete} \quad D_{KL}(P||Q) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \end{aligned}$$

Equation a.1

- KL divergence is always non-negative (as a consequence of *Gibb's Inequality*): $D_{KL}(P||Q) \geq 0$.
- When the two distributions are the same, KL divergence is zero $D_{KL}(P||P) = 0$.
- KL divergence is not symmetrical: it is easy to see that in general $D_{KL}(P||Q) \neq D_{KL}(Q||P)$, hence it is not a metric.
- Anyway a symmetric computation may be devised: $Symmetric_{KL} = D_{KL}(P||Q) + D_{KL}(Q||P)$

A.2 Score (of the log-likelihood)

It is the gradient of the log-likelihood function with respect to the parameters vector. Hence, given a certain parameter vector, the score denotes the steepness of the log-likelihood at that point in parameter space, or in other terms how much the log-likelihood would change for an infinitesimal change in the parameters.

Naming θ the parameter vector:

$$s(\theta) = \frac{\partial \log \mathcal{L}(\theta)}{\partial \theta}$$

Equation a.2

Since the log-likelihood is a function of the samples X , we can make that explicit:

$$s(\theta; x) = \frac{\partial \log \mathcal{L}(\theta; x)}{\partial \theta}$$

Equation a.3

If θ are the true parameters of the distribution $P_\theta(x)$ of samples X , the expectation of the score with respect of the distribution of X is the zero vector:

$$E_{x \sim P_\theta(x)}[s(\theta; x)] = \mathbf{0}$$

Equation a.4

To show it:

$$E_{x \sim P_\theta(x)}[s(\theta; x)] = \int_x P_\theta(x) \frac{\partial \log \mathcal{L}(\theta; x)}{\partial \theta}$$

since θ are the true parameters,
 $P_\theta(x)$ is $\mathcal{L}(\theta; x)$. Substitute it, and apply log – derivative trick (eq. 6.1)

$$\begin{aligned} &= \int_x P_\theta(x) \frac{1}{P_\theta(x)} \frac{\partial P_\theta(x)}{\partial \theta} \\ &= \int_x \frac{\partial P_\theta(x)}{\partial \theta} \end{aligned}$$

under regularity conditions it is possible to interchange the derivative and the integral
 (Leibniz Integral Rule)

$$= \frac{\partial}{\partial \theta} \int_x P_\theta(x) = \frac{\partial}{\partial \theta} 1 = 0 \quad \therefore$$

Equation a.5

A.3 Fisher Information Matrix

The Fisher Information is the variance of the score with respect to the distribution of the samples. If the parameters θ of the distribution are more than one, the Fisher Information is a matrix of covariances of the elements of the score vector, and we name it \mathbf{F}_{P_θ} .

Assuming that θ are the true parameters of the distribution of X :

$$\begin{aligned} \mathbf{F}_{P_\theta} &= \text{Var}_{x \sim P_\theta(x)}[s(\theta; x)] \\ &= E_{x \sim P_\theta(x)} \left[\left(s(\theta; x) - E_{x \sim P_\theta(x)}[s(\theta; x)] \right) \left(s(\theta; x) - E_{x \sim P_\theta(x)}[s(\theta; x)] \right)^T \right] \end{aligned}$$

from equation a.4 we know that $E_{x \sim P_\theta(x)}[s(\theta; x)] = \mathbf{0}$, hence:

$$\mathbf{F}_{P_\theta} = E_{x \sim P_\theta(x)}[s(\theta; x) s(\theta; x)^T]$$

Equation a.6

That may also be written as (recall that $s(\theta; x) = \frac{\partial \log \mathcal{L}(\theta; x)}{\partial \theta} = \nabla_\theta \log P_\theta(x)$ by definition):

$$\mathbf{F}_{P_\theta} = E_{x \sim P_\theta(x)}[\nabla_\theta \log P_\theta(x) \nabla_\theta \log P_\theta(x)^T]$$

Equation a.7

Hence, each matrix element $F_{P_\theta}[i, j]$ has the following form:

$$F_{P_\theta}[i, j] = E_{x \sim P_\theta(x)} \left[\left(\frac{\partial}{\partial \theta_i} \log P_\theta(x) \right) \left(\frac{\partial}{\partial \theta_j} \log P_\theta(x) \right) \right]$$

Equation a.8

It is clearly a square symmetric matrix.

It becomes evident that such a matrix may be estimated by samples x_i which follow the P_θ distribution :

$$\widehat{F}_{P_\theta} = \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log P_\theta(x_i) \nabla_\theta \log P_\theta(x_i)^T$$

Equation a.9

(That is actually what we do when approximating the Fisher Information Matrix for the Natural Policy Gradient in Ch.8, equation 8.3, where instead of P_θ there is the policy function π_θ).

The Fisher information matrix is always symmetric and positive semi-definite [Watanabe 2009].

A.4 Fisher Information Matrix equivalence to negative expectation of Hessian Matrix of log-probability

Now, for the next part we need to use the concept of the Hessian Matrix. Recall that the Hessian is the square matrix of second order partial derivatives of a scalar function (shortly, the Hessian is the Jacobian of the gradient). In our case the Hessian would be with respect to the parameters θ of a probability function $P_\theta(x)$, that is $\frac{\partial^2 P_\theta(x)}{\partial \theta_i \partial \theta_j}$. To simplify reading I will use the symbol $H_{P_\theta(x)}$ for it, while for the Jacobian instead I will use the symbol $\mathbf{J}()$.

If the log-likelihood function is twice differentiable with respect to θ , and under certain regularity conditions it can be shown that the Fisher Information Matrix is equivalent to the negative expected value of the Hessian matrix of the log-likelihood.

$$F_{P_\theta} = -E_{x \sim P_\theta(x)} [H_{\log P_\theta(x)}]$$

Equation a.10

So, each matrix element $F_{P_\theta}[i, j]$ is:

$$F_{P_\theta}[i, j] = -E_{x \sim P_\theta(x)} \left[\left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log P_\theta(x) \right) \right]$$

Equation a.11

To show it let us start from the Hessian of the log-likelihood, following [Kristiadi 2018]:

$$H_{\log P_\theta(x)} = \mathbf{J}(\nabla_\theta \log P_\theta(x))$$

$$= \mathbf{J} \left(\frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \right)$$

applying the quotient rule for derivatives

$$= \frac{H_{P_\theta(x)} P_\theta(x) - \nabla_\theta P_\theta(x) \nabla_\theta P_\theta(x)^T}{P_\theta(x) P_\theta(x)}$$

$$= \frac{H_{P_\theta(x)}}{P_\theta(x)} - \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \frac{\nabla_\theta P_\theta(x)^T}{P_\theta(x)}$$

let us take the expectation with respect to x distributed by $P_\theta(x)$

$$E_{x \sim P_\theta(x)}[H_{\log P_\theta(x)}] = E_{x \sim P_\theta(x)} \left[\frac{H_{P_\theta(x)}}{P_\theta(x)} - \frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \frac{\nabla_\theta P_\theta(x)^T}{P_\theta(x)} \right]$$

$$= \int_x \frac{H_{P_\theta(x)}}{P_\theta(x)} P_\theta(x) - E_{x \sim P_\theta(x)} \left[\frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \frac{\nabla_\theta P_\theta(x)^T}{P_\theta(x)} \right]$$

under regularity conditions it is possible to interchange the derivative and the integral

(Leibniz Integral Rule)

$$= H_{\int_x P_\theta(x)} - E_{x \sim P_\theta(x)} \left[\frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \frac{\nabla_\theta P_\theta(x)^T}{P_\theta(x)} \right]$$

the integral of a distribuion is = 1 and its derivative is = 0, hence $H_{\int_x P_\theta(x)} = \mathbf{0}$

$$= -E_{x \sim P_\theta(x)} \left[\frac{\nabla_\theta P_\theta(x)}{P_\theta(x)} \frac{\nabla_\theta P_\theta(x)^T}{P_\theta(x)} \right]$$

$$= -E_{x \sim P_\theta(x)} [\nabla_\theta \log P_\theta(x) \nabla_\theta \log P_\theta(x)^T]$$

$$= -\mathbf{F}_{P_\theta} \quad \therefore$$

Equation a.12

Since the Hessian Matrix is commonly used to study the curvature of a function at a critical point, the equivalence above tells us that the Fisher Information Matrix contains information about the curvature of the expectation of the log-likelihood.

A.5 Kullback-Leibler divergence approximation by second order Taylor expansion using Fisher Information Matrix

The equivalence of the Fisher Information Matrix with the negative expected value of the Hessian matrix of the log-likelihood allows to compute an approximation of the Kullback-Leibler divergence between two distributions where one is the “perturbed” version of the other (i.e. they share the same parametric form, and the parameters vector of the second distribution are obtained adding a small vector to the parameters of the first), using second order Taylor expansion for a part of the KL divergence formula.

Recall that second order Taylor series expansion for a function $f(\theta)$ is an approximation used to evaluate the function $f(\theta)$ around a certain point θ_0 , at the point $\theta_0 + \delta$:

$$f(\theta_0 + \delta) \approx f(\theta_0) + \nabla_\theta f(\theta_0)^T \delta + \frac{1}{2} \delta^T (\nabla_\theta^2 f(\theta_0)) \delta$$

Equation a.13

Consider having the distributions $P_\theta(x) = f(x; \theta)$ and $Q_{\theta, \delta}(x) = f(x; \theta + \delta)$, that means that both share the same function f , and the parameters vector of P is θ , while the parameters vector of Q is $(\theta + \delta)$, with small δ . In other terms, Q is a perturbed version of P . For instance in Reinforcement Learning P may be the old policy function and Q the new policy function obtained by an optimization step on P .

I follow the proof by [Ratliff 2013]:

$$\begin{aligned} D_{KL}(P||Q) &= D_{KL}(P_\theta(x)||Q_{\theta, \delta}(x)) = D_{KL}(f(x; \theta)||f(x; \theta + \delta)) = \\ &= \int_x f(x; \theta) \log \frac{f(x; \theta)}{f(x; \theta + \delta)} \end{aligned}$$

$$= \int_x f(x; \theta) \log f(x; \theta) - \int_x f(x; \theta) \log f(x; \theta + \delta)$$

Equation a.14

Now we apply the second order Taylor expansion only to $\log f(x; \theta + \delta)$:

$$\begin{aligned} \log f(x; \theta + \delta) &\approx \log f(x; \theta) + \nabla_{\theta} \log f(x; \theta)^T \delta + \frac{1}{2} \delta^T (\nabla_{\theta}^2 \log f(x; \theta)) \delta \\ &= \log f(x; \theta) + \left(\frac{\nabla_{\theta} f(x; \theta)}{f(x; \theta)} \right)^T \delta + \frac{1}{2} \delta^T (\nabla_{\theta}^2 \log f(x; \theta)) \delta \end{aligned}$$

Equation a.15

Now we plug a.15 into a.14:

$$\begin{aligned} D_{KL}(f(x; \theta) || f(x; \theta + \delta)) &\approx \\ &\approx \int_x f(x; \theta) \log f(x; \theta) - \int_x f(x; \theta) \left(\log f(x; \theta) + \left(\frac{\nabla_{\theta} f(x; \theta)}{f(x; \theta)} \right)^T \delta + \frac{1}{2} \delta^T (\nabla_{\theta}^2 \log f(x; \theta)) \delta \right) \\ &= \int_x f(x; \theta) \log \frac{f(x; \theta)}{f(x; \theta)} - \left(\int_x \nabla_{\theta} f(x; \theta) \right)^T \delta - \frac{1}{2} \delta^T \left(\int_x f(x; \theta) \nabla_{\theta}^2 \log f(x; \theta) \right) \delta \end{aligned}$$

Equation a.16

Now: $\int_x f(x; \theta) \log \frac{f(x; \theta)}{f(x; \theta)}$ is equal to 0 (easy to do the math or to see it as a KL divergence of two equal distributions).

Also $\left(\int_x \nabla_{\theta} f(x; \theta) \right)^T$ is equal to 0: it is possible to swap the gradient and derivative signs (under regularities) to obtain $\left(\nabla_{\theta} \int_x f(x; \theta) \right)^T$. The integral of any distribution is 1, so $\nabla_{\theta} 1 = 0$.

So it turns out:

$$D_{KL}(f(x; \theta) || f(x; \theta + \delta)) \approx -\frac{1}{2} \delta^T \left(\int_x f(x; \theta) \nabla_{\theta}^2 \log f(x; \theta) \right) \delta$$

$$\begin{aligned}
&= -\frac{1}{2} \delta^T \left(\int_x f(x; \theta) H_{\log f(x; \theta)} \right) \delta \\
&= -\frac{1}{2} \delta^T E_{x \sim f(x; \theta)} [H_{\log f(x; \theta)}] \delta
\end{aligned}$$

Equation a.17

Now, we already know that the expected value of the Hessian matrix of the log-likelihood is equal to the negative of Fisher Information Matrix (eq. a.10). So we substitute it:

$$D_{KL}(f(x; \theta) || f(x; \theta + \delta)) \approx \frac{1}{2} \delta^T \mathbf{F}_{P_\theta} \delta$$

Equation a.18

A.6 Relation between the Hessian of Kullback-Leibler divergence and Fisher Information Matrix

Again, consider having two distributions with the same function but different parameters $P_\theta(x)$ and $P_{\theta'}(x)$.

Then the gradient of Kullback-Leibler divergence between P_θ and $P_{\theta'}$ with respect to θ' would be:

$$\nabla_{\theta'} D_{KL}(P_\theta || P_{\theta'}) = \nabla_{\theta'} \int_x P_\theta(x) \log \frac{P_\theta(x)}{P_{\theta'}(x)}$$

Equation a.19

That gradient is a vector of partial derivatives with respect to all parameters of θ' . To denote the partial derivative with respect to parameter j of parametrization θ' I use the symbol $\partial_{\theta'j}$.

We have that:

$$\partial_{\theta'j} = \partial_{\theta'j} \int_x P_\theta(x) \log \frac{P_\theta(x)}{P_{\theta'}(x)}$$

under regularity conditions the derivative and integral symbol can be exchanged

$$\begin{aligned}
&= \int_x P_\theta(x) \partial_{\theta'j} \left(\log \frac{P_\theta(x)}{P_{\theta'}(x)} \right) \\
&= \int_x P_\theta(x) \frac{P_{\theta'}(x)}{P_\theta(x)} P_\theta(x) \frac{-1}{P_{\theta'}(x)^2} \partial_{\theta'j} P_{\theta'}(x)
\end{aligned}$$

$$= - \int_x \frac{P_\theta(x)}{P_{\theta'}(x)} \partial_{\theta'j} P_{\theta'}(x)$$

Equation a.20

Now let us take the derivative of eq. a.20, that is the second derivative of the KL divergence (computing at each parameter i, j): $\partial_{\theta'i} \partial_{\theta'j} d$.

$$\partial_{\theta'i} \partial_{\theta'j} d = \partial_{\theta'i} \int_x - \frac{P_\theta(x)}{P_{\theta'}(x)} \partial_{\theta'j} P_{\theta'}(x)$$

under regularity conditions the derivative and integral symbol can be exchanged

$$\begin{aligned} &= \int_x -P_\theta(x) \partial_{\theta'i} \left(\frac{\partial_{\theta'j} P_{\theta'}(x)}{P_{\theta'}(x)} \right) \\ &= \int_x -P_\theta(x) \left(\frac{\partial_{\theta'i} \partial_{\theta'j} P_{\theta'}(x) P_{\theta'}(x) - (\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)^2} \right) \\ &= \int_x P_\theta(x) \left(\frac{(\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)^2} - \frac{\partial_{\theta'i} \partial_{\theta'j} P_{\theta'}(x)}{P_{\theta'}(x)} \right) \end{aligned}$$

now, if we evaluate that integral at the parameter point $\theta' = \theta$ we obtain :

$$\begin{aligned} \partial_{\theta'i} \partial_{\theta'j} d \Big|_{\theta' = \theta} &= \int_x P_\theta(x) \left(\frac{(\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)^2} - \frac{\partial_{\theta'i} \partial_{\theta'j} P_{\theta'}(x)}{P_{\theta'}(x)} \right) \Big|_{\theta' = \theta} \\ &= \int_x \frac{(\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)} - \partial_{\theta'i} \partial_{\theta'j} P_{\theta'}(x) \Big|_{\theta' = \theta} \end{aligned}$$

it is easy to see that $\int_x \partial_{\theta'i} \partial_{\theta'j} P_{\theta'}(x) = \partial_{\theta'i} \partial_{\theta'j} \int_x P_{\theta'}(x) = \partial_{\theta'i} \partial_{\theta'j} 1 = 0$

$$\begin{aligned} &= \int_x \frac{(\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)} \Big|_{\theta' = \theta} \\ &= \int_x \frac{P_{\theta'}(x)}{P_{\theta'}(x)} \frac{(\partial_{\theta'i} P_{\theta'}(x)) (\partial_{\theta'j} P_{\theta'}(x))}{P_{\theta'}(x)} \Big|_{\theta' = \theta} \end{aligned}$$

$$\begin{aligned}
&= \int_x P_{\theta'}(x) (\partial_{\theta'_i} \log P_{\theta'}(x)) (\partial_{\theta'_j} \log P_{\theta'}(x)) \Big|_{\theta' = \theta} \\
&= E_{x \sim P_{\theta}(x)} [(\partial_{\theta_i} \log P_{\theta}(x)) (\partial_{\theta_j} \log P_{\theta}(x))] \\
&= \mathbf{F}_{P_{\theta}}[i, j] \quad (\text{see eq. a.8}) \\
&\Rightarrow H_{\theta' D_{KL}(P_{\theta} || P_{\theta'})} \Big|_{\theta' = \theta} = \mathbf{F}_{P_{\theta}}
\end{aligned}$$

Equation a.21

That means that the Hessian of Kullback-Leibler divergence between two distributions of the same family $P_{\theta}(x)$ and $P_{\theta'}(x)$, with respect to θ' , evaluated at $\theta' = \theta$ is equal to the Fisher Information Matrix of $P_{\theta}(x)$.

References

- [Amari 1998] Amari, S. "Natural gradient works efficiently in learning." Neural Computation, 10, 251. (1998).
- [Bagnell and Schneider 2003] Bagnell, J., & Schneider, J. "Covariant policy search." In Proceedings of the international joint conference on artificial intelligence (pp. 1019–1024). (2003).
- [Baird 1993] Baird, L. C. "Advantage Updating." Wright Lab. 1993, Technical Report WL-TR-93-1146.
- [Degris et al. 2012] Degris, T., White, M., and Sutton, R. S. "Linear off-policy actor-critic." In 29th International Conference on Machine Learning.
- [Glorot and Bengio 2010] Glorot, X., & Bengio, Y. "Understanding the difficulty of training deep feedforward neural networks. ", In Proceedings of the thirteenth international conference on artificial intelligence and statistics (pp. 249-256). JMLR Workshop and Conference Proceedings. (2010)
- [Irpan 2018] Alex Irpan, "Deep Reinforcement Learning Doesn't Work Yet", (2018) <https://www.alexirpan.com/2018/02/14/rl-hard.html>
- [Kakade 2002] Kakade Sham, "A natural policy gradient", Advances in Neural Information Processing Systems, pp. 1057–1063. MIT Press, 2002.

[Kakade and Langford 2002] Kakade, Sham and Langford, John. Approximately optimal approximate reinforcement learning. In ICML, volume 2, pp. 267–274, 2002.

[Kristiadi 2018] Kristiadi Augustinus,
<https://agustinus.kristia.de/techblog/2018/03/11/fisher-information/>

[Levine 2021a] Sergey Levine, CS 285 at UC Berkeley, Deep Reinforcement Learning, Lecture 9 (part 4)
<http://rail.eecs.berkeley.edu/deeprlcourse-fa21/static/slides/lec-9.pdf>
<https://www.youtube.com/watch?v=QWnpF0FaKL4>

[Levine 2021b] Sergey Levine, CS 285 at UC Berkeley, Deep Reinforcement Learning, Lecture 15 (part 2)
<https://rail.eecs.berkeley.edu/deeprlcourse/static/slides/lec-15.pdf>
<https://www.youtube.com/watch?v=9HrN6nHoxD8>

[Levine and Koltun 2013] Sergey Levine and Vladlen Koltun, "Guided Policy Search", ICML 2013, PMLR 28(3):1-9, 2013.
<http://proceedings.mlr.press/v28/levine13.pdf>

[Mnih et al. 2013] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. , "Playing Atari with Deep Reinforcement Learning." (2013), arXiv:1312.5602.

[OpenAi 2018A] OpenAi
Vanilla policy gradient and Expected Grad-Log-Prog Lemma:
https://spinningup.openai.com/en/latest/spinningup/rl_intro3.html

[OpenAi 2018B] OpenAi
Reward-to-go proof
https://spinningup.openai.com/en/latest/spinningup/extra_pg_proof1.html

[OpenAi 2018C] OpenAi
Q-function in policy gradient proof:
https://spinningup.openai.com/en/latest/spinningup/extra_pg_proof2.html

[Peters 2007] Peters, J. "Machine learning of motor skills for robotics. "Ph.D. thesis University of Southern California, Los Angeles, CA, 90089, USA. (2007).

[Peters and Schaal 2008] Jan Peters, Stefan Schaal , "Reinforcement learning of motor skills with policy gradients." , Neural networks, 21(4), 682-697. (2008)

[Peters et al. 2003] Peters, J., Vijayakumar, S., & Schaal, S. "Reinforcement learning for humanoid robotics." In Proceedings of the IEEE-RAS international conference on humanoid robots (HUMANOIDS) (pp. 103–123).(2003)

[Ratliff 2013] Ratliff, Nathan. "Information Geometry and Natural Gradients" in Mathematics for Intelligent Systems.
<https://www.nathanratliff.com/pedagogy/mathematics-for-intelligent-systems>
<https://drive.google.com/file/d/1jDM9yZI1KrtH3JJznPE7S1QzNZZB2xs2/view>
https://web.archive.org/web/20200924055016/https://ipvs.informatik.uni-stuttgart.de/mlr/wp-content/uploads/2015/01/mathematics_for_intelligent_systems_lecture12_notes_l.pdf

[Saxe et al. 2013] Saxe, A. M., McClelland, J. L., & Ganguli, S. , "Exact solutions to the nonlinear dynamics of learning in deep linear neural networks." ICLR (2013), arXiv preprint arXiv:1312.6120.

[Schulman et al. 2015] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, Pieter Abbeel , "Trust Region Policy Optimization", In International conference on machine learning (pp. 1889-1897). PMLR, arXiv:1502.05477

[Schulman et al. 2016] John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan and Pieter Abbeel, "High-Dimensional Continuous Control Using Generalized Advantage Estimation", ICLR 2016, arXiv:1506.02438

[Schulman et al. 2017] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). , "Proximal policy optimization algorithms" , arXiv preprint arXiv:1707.06347.

[Silver et al. 2014] Silver D., Lever G., Heess N., Degris T., Wierstra D., & Riedmiller M. "Deterministic policy gradient algorithms". In International conference on machine learning (pp. 387-395). PMLR. (2014, January)

[Soemers 2019] Dennis Soemers, "Reward-to-go proof"
<https://ai.stackexchange.com/questions/9614/why-does-the-reward-to-go-trick-in-policy-gradient-methods-work/10369>

[Sutton & Barto 2018] Richard C. Sutton, Andrew C. Barto, "Reinforcement Learning - An introduction. 2nd Ed.", MIT Press - ISBN: 9780262039246
<http://incompleteideas.net/book/RLbook2020.pdf>

[Sutton et al. 1999] Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. (1999). "Policy gradient methods for reinforcement learning with function approximation." In Neural Information Processing Systems 12, pages 1057–1063.

[Szepesvári 2010] Csaba Szepesvári, "Algorithms for Reinforcement Learning", Morgan & Claypool 2010, ISBN: 9781608454921

[van Hasselt 2010] Hado van Hasselt, "Double Q-learning."
Advances in neural information processing systems, 23, 2613-2621. (2010)

[Watanabe 2009] Watanabe Sumio, "Algebraic Geometry and Statistical Learning Theory", Cambridge University Press, 2009. ISBN-13 978-0-521-86467-1

[Watkins 1989] Christopher J. C. H. Watkins, "Learning from delayed rewards.", PhD Thesis, King's College (1989).

[Williams 1992] Williams, R. J. , "Simple statistical gradient-following algorithms for connectionist reinforcement learning.", Machine learning, 8(3), 229-256. (1992)