

POLICY GRADIENT ALGORITHMS FROM PROBABILITY THEORY PERSPECTIVE

JULIAN WERGIELUK

ABSTRACT. This short note provides a concise description of the model architecture and learning algorithms of the agent developed in this project. We also report learning performance of the agent and provide a list of possible future model improvements.

1. DESCRIPTION OF THE LEARNING ALGORITHM

1.1. Summary of the notation. In this write-up I am going to cast the core ideas of policy gradient methods into a rigorous probabilistic framework. Much of the confusion and difficulty of understanding reinforcement learning comes from sloppy and incomplete mathematical notation.

Let's setup up the stage for the show and define a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. Consider a Markov Decision Process (MDP) with measurable state and action spaces $(\mathbb{S}, \mathcal{S})$ and $(\mathbb{A}, \mathcal{A})$. A policy π is a Markov kernel from $(\mathbb{S}, \mathcal{S})$ to $(\mathbb{A}, \mathcal{A})$, i.e. for each $s \in \mathbb{S}$, $\pi(\cdot | s)$ is a probability distribution on \mathcal{A} . A trajectory

$$\tau = (S_0, A_0, R_1, S_1, A_1, R_2, \dots)$$

encodes a sequence of states, actions and resulting numerical rewards pertaining to a policy π . The elements of τ are random variables, such that S_0 has the *start-state distribution* ρ_0 (which does not depend on the policy π), $\pi(\cdot | s)$ is the conditional distribution of A_i given $S_i = s$ for any $s \in \mathbb{S}$, and $\mathbb{P}(\cdot | S_{i-1} = s, A_{i-1} = a)$ is the conditional distribution of S_i . The rewards R_i are random variables taking values in \mathbb{R} . The return or discounted future reward G_t at $i \in \mathbb{N}$ is defined as

$$(1) \quad G_i = \sum_{k=0}^{\infty} \gamma^k R_{i+k+1},$$

where γ is a hyperparameter from the interval $(0, 1)$, used to guarantee the convergence of the series (1).

Since τ is a sequence of random variables, and therefore itself a random variable taking values in the trajectory space

$$\mathbb{T} = \mathbb{S} \times \mathbb{A} \times (\mathbb{S} \times \mathbb{A} \times \mathbb{R})^{\infty},$$

we can sample from τ to obtain sequences of the form

$$(s_0, a_0, r_1, s_1, a_1, r_2, \dots) \in \mathbb{T}.$$

These samples of τ are called *episodes* or *rollouts*.

The central object of interest is the mechanics of the world encoded in the transition probability measure $\mathbb{P}(\cdot | s, a)$. For $i > 0$, it is the distribution of S_i under the conditions $S_{i-1} = s$ and $A_{i-1} = a$ for fixed $s \in \mathbb{S}$ and $a \in \mathbb{A}$. This distribution does not depend on i , and, more importantly, does not depend on the policy π .

To simplify the setup, we assume that the reward R_i at time i is a deterministic function of the relevant states and actions. Depending on the problem at hand these are the three common choices encountered in the literature:

$$\begin{aligned} R_i &= \phi_1(S_i), \\ R_i &= \phi_2(A_{i-1}, S_i), \\ R_i &= \phi_3(S_{i-1}, A_{i-1}, S_i) \end{aligned}$$

ϕ_1 is probably the most frequently used form.

In reinforcement learning, we parametrize the policy π_θ with an \mathbb{R}^d vector θ , $d \geq 1$. The goal is to find good values for θ which lead to high cumulative expected reward

$$J(\theta) = \mathbb{E}_\theta[G_0] = \mathbb{E}_\theta \left[\sum_{k=0}^{\infty} \gamma^k R_{k+1} \right].$$

The subscript θ next to the expectation operator indicates that the underlying probability measure \mathbb{P}_θ combines both \mathbb{P} and π_θ .

1.2. Value functions. The *state-value function* is defined for each state $s \in \mathbb{S}$ as

$$(2) \quad V_\theta(s) = \mathbb{E}_\theta[G_0 | S_0 = s] = \mathbb{E}_\theta \left[\sum_{k=0}^{\infty} \gamma^k R_{k+1} | S_0 = s \right].$$

The *action-value function* is defined for each state-action pair $(s, a) \in \mathbb{S} \times \mathbb{A}$ as

$$(3) \quad Q_\theta(s, a) = \mathbb{E}_\theta[G_0 | S_0 = s, A_0 = a] = \mathbb{E}_\theta \left[\sum_{k=0}^{\infty} \gamma^k R_{k+1} | S_0 = s, A_0 = a \right].$$

The *advantage function* is the difference

$$(4) \quad A_\theta(s, a) = V_\theta(s) - Q_\theta(s, a).$$

1.3. Episodic and perpetual tasks. Reinforcement learning tasks can be partitioned into two categories, depending on whether the task can be completed. A task is called *episodic* if there is an end of the task, i.e. some notion of “being done”. If there is no

such notion, the task is called *perpetual*¹. Typical examples of episodic tasks are games. Most board games like chess and go, and computer games fall into this category.

It is important to note, that many episodic tasks can also run forever. For example, in chess, we can design a policy that “runs in circles” and does not pursue the goal of winning the game. Letting such a policy play against itself will result in an endless game.

What distinguishes an episodic task from a perpetual one, is the existence of a non-empty set of end states Δ in \mathbb{S} , aptly called the *coffin space*, for which we require that $\mathbb{P}(S_{i+1} = S_i | S_i \in \Delta) = 1$ and $\mathbb{P}(R_i = 0 | S_i \in \Delta) = 1$. This ensures that all the states in Δ are absorbing and that no additional reward can be accumulated after the task has been completed.

1.4. Categorical and continuous action spaces. The cardinality of the action space has a huge impact on the development of the theory of policy gradients as well as on the design and implementation of the policy optimization algorithms. Essentially, there are two cases we need to deal with. First, if the action space has finite number of elements, we call it *categorical*. Finite number of possible actions implies that the policy $\pi(.|s)$ evaluated at the state s , can be identified with a vector with $|\mathbb{A}|$ number of elements. Also, it is convenient to define the likelihood function f of the policy π as

$$f(a|s) = \pi(\{a\}|s), a \in \mathbb{A}.$$

$f(.|s)$ is the probability function of the probability measure $\pi(.|s)$.

The second case encountered in the literature is the *continuous* action space \mathbb{A} which is a “nice” subset of $\mathbb{R}^{n_{\mathbb{A}}}$. Typical example of a continuous action space is a product of finite or infinite intervals in $\mathbb{R}^{n_{\mathbb{A}}}$. Let us assume that, for each state s , the measure $\pi(.|s)$ has a density $f(.|s)$, i.e. for an action set $\eta \in \mathcal{A}$ we have

$$\pi(\eta|s) = \int_{\eta} f(a|s) da.$$

As in the categorical case, we call $f(.|s)$ the likelihood function of the measure $\pi(.|s)$.

1.5. Technical assumptions. In the following list, we gather the technical assumptions.

Assumption 1.1. (1) *The reward function ϕ is measurable and bounded.*

¹In the literature, people use the adjective “continuous”. But we already have “continuous spaces”, “continuous time”, “continuous function”, etc. The word “perpetual” is much better at conveying the property of running indefinitely long. See e.g. <https://www.youtube.com/watch?v=t5TK1OxEj9A>.

Let us state that our body of thought compiles without errors.

Theorem 1.2. *Suppose the conditions stated in Assumption 1.1 hold. Then we have:*

(1) *The value functions V_θ , Q_θ , and A_θ are well-defined.*

1.6. Cumulative reward optimization. One of our goals is to find the values of the parameter θ yielding a policy that consistently produces high cumulative expected award $J(\theta)$.

The gradient of $J(\theta)$ is given by

$$\nabla_\theta J(\theta) = \nabla_\theta \mathbb{E}_\theta [G_0] = \mathbb{E}_\theta \left[G_0 \sum_{i=0}^{\infty} \nabla_\theta \log \pi_\theta(A_i | S_i) \right].$$

Using the Markov property we get

$$\nabla_\theta J(\theta) = \mathbb{E}_\theta \left[\sum_{i=0}^{\infty} G_i \nabla_\theta \log \pi_\theta(A_i | S_i) \right].$$

Policy gradient with the state-value baseline

$$\nabla_\theta J(\theta) = \mathbb{E}_\theta \left[\sum_{i=0}^{\infty} (G_i - V_\theta(S_i)) \nabla_\theta \log \pi_\theta(A_i | S_i) \right].$$

Policy gradient with the advantage function baseline

$$\nabla_\theta J(\theta) = \mathbb{E}_\theta \left[\sum_{i=0}^{\infty} A_\theta(S_i, A_i) \nabla_\theta \log \pi_\theta(A_i | S_i) \right].$$

2. TRAINING ANALYSIS

3. IDEAS FOR FUTURE WORK

REFERENCES

- [1] Erhan Çinlar. *Probability and Stochastics*. 1st ed. Springer Verlag, 2011, p. 572. ISBN: 0387878580.

Email address: julian.wergieluk@risklab.com