

Why does the "reward to go" trick in policy gradient methods work?

Asked 6 years, 3 months ago Modified 1 year, 3 months ago Viewed 7k times



In the policy gradient method, there's a trick to reduce the variance of policy gradient. We use causality, and remove part of the sum over rewards so that only actions happened after the reward are taken into account (See here



http://rail.eecs.berkeley.edu/deeprlcourse/static/slides/lec-5.pdf, slide 18).



Why does it work? I understand the intuitive explanation, but what's the rigorous proof of it? Can you point me to some papers?



reinforcement-learning math policy-gradients rewards reward-to-go

Share Improve this question Follow

edited Oct 10, 2020 at 15:52



asked Dec 20, 2018 at 1:00

Konstantin Solomatov

3 Answers

Sorted by: Highest score (default)



An important thing we're going to need is what is called the "Expected Grad-Log-Prob Lemma" <u>here</u>" (proof included on that page), which says that (for any t):



$$\mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t)
ight] = 0.$$



Taking the analytical expression of the gradient (from, for example, slide 9) as a starting point:







$$egin{aligned}
abla_{ heta} J(heta) &= \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[\left(\sum_{t=1}^{T}
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \right) \left(\sum_{t=1}^{T} r(s_t, a_t)
ight)
ight] \ &= \sum_{t=1}^{T} \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=1}^{T} r(s_{t'}, a_{t'})
ight] \ &= \sum_{t=1}^{T} \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=1}^{T} r(s_{t'}, a_{t'}) +
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=t}^{T} r(s_{t'}, a_{t'})
ight] \ &= \sum_{t=1}^{T} \left(\mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=t}^{t-1} r(s_{t'}, a_{t'})
ight] \ &+ \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=t}^{T} r(s_{t'}, a_{t'})
ight]
ight) \end{aligned}$$

At the t^{th} "iteration" of the outer sum, the random variables $\sum_{t'=1}^{t-1} r(s_{t'}, a_{t'})$ and $\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t)$ are independent (we assume, by definition, the action only depends on the most recent state), which means we are allowed to split the expectation:

$$egin{aligned}
abla_{ heta} J(heta) &= \sum_{t=1}^T \left(\mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[\sum_{t'=1}^{t-1} r(s_{t'}, a_{t'})
ight] \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t)
ight] \\ &+ \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=t}^T r(s_{t'}, a_{t'})
ight]
ight) \end{aligned}$$

The first expectation can now be replaced by 0 due to the lemma mentioned at the top of the post:

$$egin{aligned}
abla_{ heta} J(heta) &= \sum_{t=1}^T \mathbb{E}_{ au \sim \pi_{ heta}(au)} \left[
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \sum_{t'=t}^T r(s_{t'}, a_{t'})
ight] \ &= \mathbb{E}_{ au \sim \pi_{ heta}(au)} \sum_{t=1}^T
abla_{ heta} \log \pi_{ heta}(a_t \mid s_t) \left(\sum_{t'=t}^T r(s_{t'}, a_{t'})
ight). \end{aligned}$$

The expression on slide 18 of the linked slides is an unbiased, sample-based estimator of this gradient:

$$abla_{ heta} J(heta) pprox rac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T}
abla_{ heta} \log \pi_{ heta}(a_{i,t} \mid s_{i,t}) \left(\sum_{t'=t}^{T} r(s_{i,t'}, a_{i,t'})
ight)$$

For a more formal treatment of the claim that we can pull $\sum_{t'=1}^{t-1} r(s_{t'}, a_{t'})$ out of an expectation due to the Markov property, see this page:

https://spinningup.openai.com/en/latest/spinningup/extra_pg_proof1.html

Share Improve this answer Follow

edited Nov 25, 2020 at 20:14

Mike Land

103 3

answered Feb 2, 2019 at 19:16

Dennis Soemers ♦

1 Why is the form of "Expected Grad-Log-Prob Lemma" mentioned at the top of the post legit? It is different than the form of the proof shown in the link (the proof in the link has no problem). Your lemma is an expectation over **trajectories**. – starriet 차주녕 Jan 4, 2022 at 3:09 ▶

@starriet For the individual choice of a single action at one point in time within a trajectory, the trajectory as a whole has no influence (our policy does not look at the history or the future, just at the current state), so $\mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \right] = \mathbb{E}_{a_t \sim \pi_{\theta}} \left[\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \right]$, where the right-hand side more closely matches that notation from the link. – Dennis Soemers \bullet Jan 4, 2022 at 9:43 \r

@starriet However, we use the trajectory-based notation because in the larger equation, we have these sums over all time steps, and that's the part where the full trajectory actually becomes relevant.
Dennis Soemers ♦ Jan 4, 2022 at 9:45

Thanks, I think I got your point. If I understood correctly, the right-hand side is at a *specific state*, and the left-hand side is at the various states and actions of *all possible trajectories*, right? Since the RHS is 0 and it's about one state, LHS is also 0 because LHS is an expectation over all states (this is how I understood. please let me know if I'm wrong). – starriet 차주녕 Jan 4, 2022 at 12:41

But still, It's difficult to understand why the form of the LHS is legit. It's an expectation over trajectories, but the term inside of the square bracket is a function of action and state. Given that a trajectory is composed of many states and actions, how the expectation over trajectories is possible if the term inside of the brackets is smaller than a trajectory? – starriet 차주녕 Jan 4, 2022 at 12:51



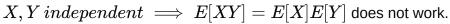
1

A correct proof is given at





On the contrary, an attempt to prove "reward to go" by applying





The random variables $\sum_{t'=1}^{t-1} r(s_{t'}, a_{t'})$ and $\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t)$ are generally not independent. They would be independent, if

$$P[\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) \mid \sum_{t'=1}^{t-1} r(s_{t'}, a_{t'})] = P[\nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t)] \text{ but this does not follow from } \pi_{\theta}(a_t \mid s_t) = \pi_{\theta}(a_t \mid s_t, s_{t-1}, \dots, s_1, a_{t-1}, \dots, a_1)$$

Consider the following setup as an example for the dependency between the two random variables in question.

$$\pi(a_t = 0 \mid s_t) = \frac{1}{s_t}$$

 $r(s_t,a_t)=1$ i.e. the event reward at time t equals 1 and the total reward from time 1 to time t-1 equals t-1

 $P(s_1=1)=1$ i.e. all trajectories start in state 1

 $P(s_{t+1}=s_t+1\mid s_t,a_t)=1$ i.e. the next state, given the current state and the current action, is always the current state + 1

As a result, the higher time t, the higher the state s_t , the higher the probability for action $a_t=1$ and the higher the total reward from time 0 to time t-1, i.e. the values of $\pi(a_t \mid s_t)$ and thus the values of $\nabla_{\theta} \log \pi(a_t \mid s_t)$ are highly correlated with the total rewards from time 0 to time t-1, hence respective random variables are not independent of each other.

Share Improve this answer Follow

answered Sep 27, 2022 at 14:36





By linearity of expection and using the distributive property:

$$egin{aligned}
abla_{m{ heta}} J(m{ heta}) &= \mathbb{E}_{ au\sim p_{m{ heta}}(au)} \left[\left(\sum_{t=1}^{T}
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)
ight)
ight] \ &= \sum_{t=1}^{T} \mathbb{E}_{ au\sim p_{m{ heta}}(au)} \left[
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \sum_{t'=1}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})
ight] \ &= \sum_{t=1}^{T} \mathbb{E}_{ au\sim p_{m{ heta}}(au)} \left[
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \sum_{t'=1}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) +
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})
ight] + \ &\sum_{t=1}^{T} \mathbb{E}_{ au\sim p_{m{ heta}}(au)} \left[
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})
ight] \ \end{pmatrix}$$

and the first term can be further expanded as follows:

$$\begin{split} &\sum_{t=1}^{T} \sum_{t'=1}^{t-1} \mathbb{E}_{\tau \sim p_{\boldsymbol{\theta}}(\tau)} \left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{t} \mid \mathbf{s}_{t}) \, r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \\ &= \sum_{t=1}^{T} \sum_{t'=1}^{t-1} \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}, \mathbf{s}_{t'}, \mathbf{a}_{t'}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_{t}, \mathbf{a}_{t}, \mathbf{s}_{t'}, \mathbf{a}_{t'})} \left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{t} \mid \mathbf{s}_{t}) \, r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \\ &= \sum_{t=1}^{T} \sum_{t'=1}^{t-1} \mathbb{E}_{(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_{t'}, \mathbf{a}_{t'})} \left[\mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_{t}, \mathbf{a}_{t} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'})} \left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{t} \mid \mathbf{s}_{t}) \, r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \\ &= \sum_{t=1}^{T} \sum_{t'=1}^{t-1} \mathbb{E}_{(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_{t'}, \mathbf{a}_{t'})} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \, \mathbb{E}_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim p_{\boldsymbol{\theta}}(\mathbf{s}_{t}, \mathbf{a}_{t} \mid \mathbf{s}_{t'})} \left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{t} \mid \mathbf{s}_{t}) \right] \right] \end{split}$$

The innermost expectation can be broken down further into

$$egin{aligned} & \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim p_{oldsymbol{ heta}(\mathbf{s}_t, \mathbf{a}_t \mid \mathbf{s}_{t'}, \mathbf{a}_{t'})} \left[
abla_{oldsymbol{ heta}} \log \pi_{oldsymbol{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t)
ight] \ &= \int \int
abla_{oldsymbol{ heta}} \log \pi_{oldsymbol{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \pi_{oldsymbol{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) p(\mathbf{s}_t \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) d\mathbf{a}_t \ d\mathbf{s}_t \ &= \int p(\mathbf{s}_t \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) \int
abla_{oldsymbol{ heta}} \log \pi_{oldsymbol{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \pi_{oldsymbol{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) d\mathbf{a}_t \ d\mathbf{s}_t \end{aligned}$$

where $\int \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \, \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \, d\mathbf{a}_t = \mathbb{E}_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t)} \left[\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t) \right] = 0$ due to the EGLP lemma

As a result, we have

$$abla_{m{ heta}} J(m{ heta}) = \sum_{t=1}^T \mathbb{E}_{ au \sim p_{m{ heta}}(au)} \left[
abla_{m{ heta}} \log \pi_{m{ heta}}(\mathbf{a}_t \mid \mathbf{s}_t) \, \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})
ight]$$

Share Improve this answer Follow

edited Jan 5, 2024 at 15:52

answered Jan 5, 2024 at 15:45



Jus Forever

Start asking to get answers

Find the answer to your question by asking.

Ask question

Explore related questions

reinforcement-learning math policy-gradients rewards reward-to-go

See similar questions with these tags.