

CS 234: Assignment #3

Winter 2019

Li Quan Khoo (SCPD)

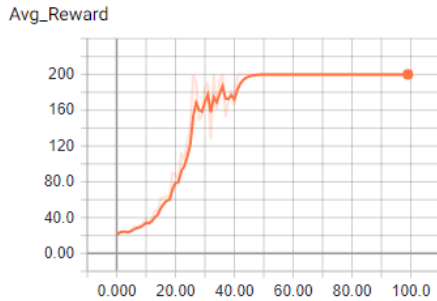
1 Policy Gradient Methods (50 pts + 15 pts writeup)

1.1 Coding Questions (50 pts)

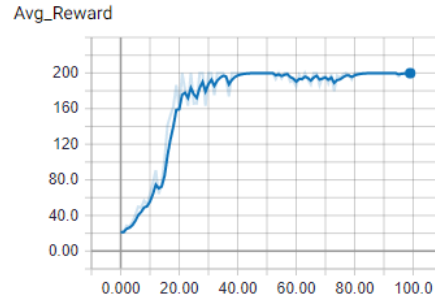
- <https://github.com/lqkhoo/cs234-winter-2019/blob/master/assignment3/pg.py>

1.2 Writeup Questions (15 pts)

(a) (4 pts) (CartPole-v0)



(a) CartPole-v0 with baseline. Average reward over 100 iterations.



(b) CartPole-v0 without baseline. Average reward over 100 iterations.

Where

$$r(\tau) = \sum_t r(s_t, a_t), \quad \text{and} \quad J(\theta) = \mathbb{E}_{\tau \sim P_\theta(\tau)}[r(\tau)],$$

We have derived the Monte-Carlo estimate of the gradient of $J(\theta)$ for policy gradient to be:

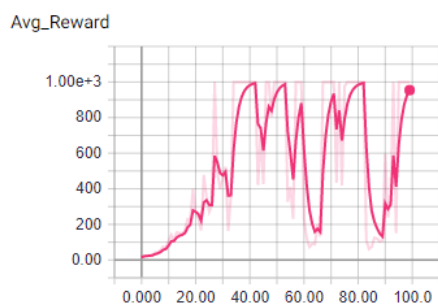
$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \pi_\theta(\tau) [r(\tau) - b]$$

Where b is some constant scalar bias of the reward function. Mathematically, the baseline is just another bias, and we already know that this gradient $\nabla_\theta J(\theta)$ is indifferent to the bias *in expectation*. However, the actual size and value of the bias still influences convergence (consider a bias close to $+\infty$ or $-\infty$) due to how the probability mass of the policy is shifted when we perform the gradient update, as well as the variance of the quantity $r(\tau) - b$. (Berkeley CS294, lecture 4)

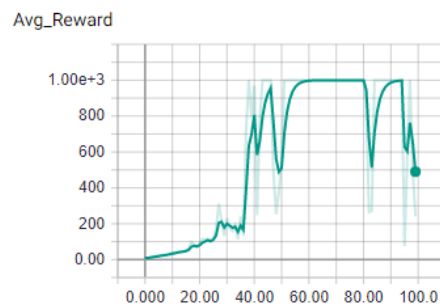
Here, the baseline we are taking is the state value function $V(s_t)$, which represents the average of returns of all actions from that state. Suppose we have a good estimate of $V(s_t)$. Mathematically, $V(s_t)$ already summarizes all future returns, and the advantage $A(s_t, a_t) = r(s_t, a_t) - V(s_t)$ becomes insulated from the variance of returns of all future time steps *in the current trajectory*. This lowers the variance of our estimated action values. Semantically, we are shifting the probability mass of the policy to make actions which are *better than average* (our chosen baseline) to be more likely.

Empirically, we can see from the plots above that using the baseline tends to give "better" convergence, although this is far from being guaranteed. Our observation here is that the plateau is more stable. The standard deviation of rewards (not shown here) tends to be lower as well.

(b) (4 pts) (InvertedPendulum-v1)



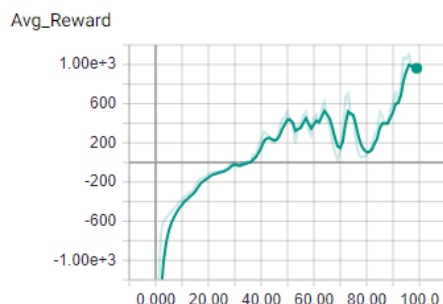
(a) Pendulum-v1 with baseline. Average reward over 100 iterations.



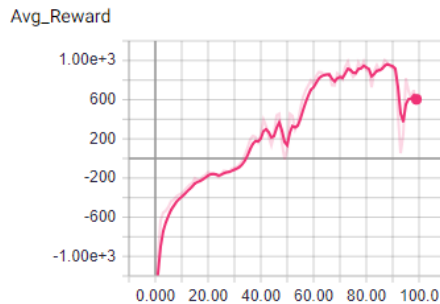
(b) Pendulum-v1 without baseline. Average reward over 100 iterations.

In this case, we notice that running with baseline tends to have faster convergence, but here, it doesn't seem to help very much with the instability of the algorithm. This instability is due to the high variance of REINFORCE (Monte-Carlo "vanilla" policy gradient), relative to comparable methods such as natural policy gradient and TRPO.

(c) (7 pts) (HalfCheetah-v1)



(a) HalfCheetah-v1 with baseline. Average reward over 100 iterations.



(b) HalfCheetah-v1 without baseline. Average reward over 100 iterations.

In this particular run, the baseline didn't seem to help, although on average it's supposed to have the same properties as explained in (a) by reducing the variance of the gradient.

2 Best Arm Identification in Multiarmed Bandit (35pts)

In this problem we focus on the Bandit setting with rewards bounded in $[0, 1]$. A Bandit problem instance is defined as an MDP with just one state and action set \mathcal{A} . Since there is only one state, a “policy” consists of the choice of a single action: there are exactly $A = |\mathcal{A}|$ different deterministic policies. Your goal is to design a simple algorithm to identify a near-optimal arm with high probability.

Imagine we have n samples of a random variable x , $\{x_1, \dots, x_n\}$. We recall Hoeffding’s inequality below, where \bar{x} is the expected value of a random variable x , $\hat{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is the sample mean (under the assumption that the random variables are in the interval $[0, 1]$), n is the number of samples and $\delta > 0$ is a scalar:

$$\Pr \left(|\hat{x} - \bar{x}| > \sqrt{\frac{\log(2/\delta)}{2n}} \right) < \delta.$$

Assuming that the rewards are bounded in $[0, 1]$, we propose this simple strategy: allocate an identical number of samples $n_1 = n_2 = \dots = n_A = n_{des}$ to every action, compute the average reward (empirical payout) of each arm $\hat{r}_{a_1}, \dots, \hat{r}_{a_A}$ and return the action with the highest empirical payout $\arg \max_a \hat{r}_a$. The purpose of this exercise is to study the number of samples required to output an arm that is at least ϵ -optimal with high probability. Intuitively, as n_{des} increases the empirical payout \hat{r}_a converges to its expected value \bar{r}_a for every action a , and so choosing the arm with the highest empirical payout \hat{r}_a corresponds to approximately choosing the arm with the highest expected payout \bar{r}_a .

- (a) (15 pts) We start by defining a *good event*. Under this *good event*, the empirical payout of each arm is not too far from its expected value. Starting from Hoeffding inequality with n_{des} samples allocated to every action show that:

$$\Pr \left(\exists a \in \mathcal{A} \quad s.t. \quad |\hat{r}_a - \bar{r}_a| > \sqrt{\frac{\log(2/\delta)}{2n_{des}}} \right) < A\delta.$$

In other words, the *bad event* is that at least one arm has an empirical mean that differs significantly from its expected value and this has probability at most $A\delta$.

- (b) (20 pts) After pulling each arm (action) n_{des} times our algorithm returns the arm with the highest empirical payout:

$$a^\dagger = \arg \max_a \hat{r}_a$$

Notice that a^\dagger is a random variable. Define a^* as the optimal arm (that yields the highest average reward $a^* = \arg \max_a \bar{r}_a$). Suppose that we want our algorithm to return at least an ϵ optimal arm with probability $1 - \delta'$, as follows:

$$\Pr \left(\bar{r}_{a^\dagger} \geq \bar{r}_{a^*} - \epsilon \right) \geq 1 - \delta'.$$

How many samples are needed to ensure this? Express your result as a function of the number of actions A , the required precision ϵ and the failure probability δ' .