

[SOLUTION TEMPLATE] Assignment 2: Policy Gradients

Author: Jeremiah Chen

4 Policy Gradients

- Create two graphs:
 - In the first graph, compare the learning curves (average return vs. number of environment steps) for the experiments prefixed with `cartpole`. (The small batch experiments.)
 - In the second graph, compare the learning curves for the experiments prefixed with `cartpole_lb`. (The large batch experiments.)

For all plots in this assignment, the x -axis should be number of environment steps, logged as `Train_EnvstepsSoFar` (*not* number of policy gradient iterations).

- Answer the following questions briefly:
 - Which value estimator has better performance without advantage normalization: the trajectory-centric one, or the one using reward-to-go?
 - Did advantage normalization help?
 - Did the batch size make an impact?
- Provide the exact command line configurations (or `#@params` settings in Colab) you used to run your experiments, including any parameters changed from their defaults.

q2_pg_cartpole_CartPole-v0_06-06-2025_00-19-16	●
q2_pg_cartpole_rtg_CartPole-v0_06-06-2025_00-36-28	●
q2_pg_cartpole_lb_CartPole-v0_06-06-2025_00-49-55	●
q2_pg_cartpole_lb_rtg_CartPole-v0_06-06-2025_00-51-51	●
q2_pg_cartpole_na_CartPole-v0_06-06-2025_01-02-58	●
q2_pg_cartpole_rtg_na_CartPole-v0_06-06-2025_01-03-42	●
q2_pg_cartpole_lb_na_CartPole-v0_06-06-2025_01-04-15	●
q2_pg_cartpole_lb_rtg_na_CartPole-v0_06-06-2025_01-06-09	●

Figure 1: Labels for CartPole experiments

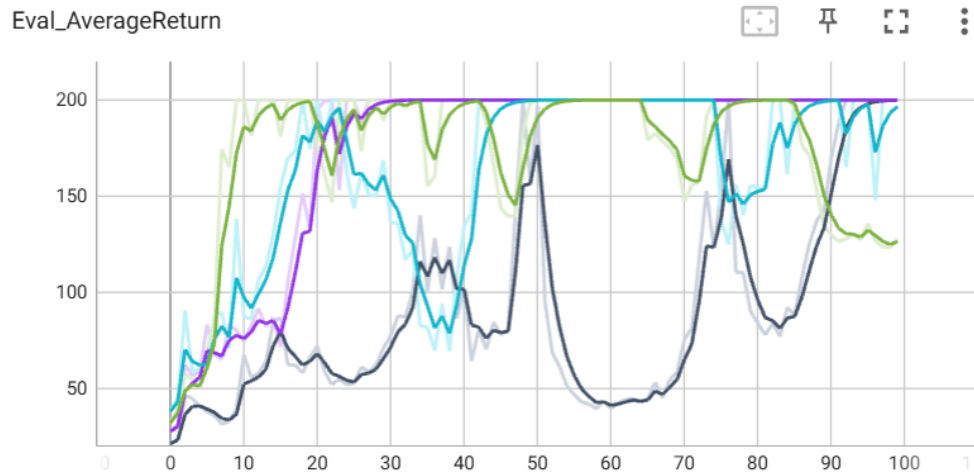


Figure 2: CartPole experiments

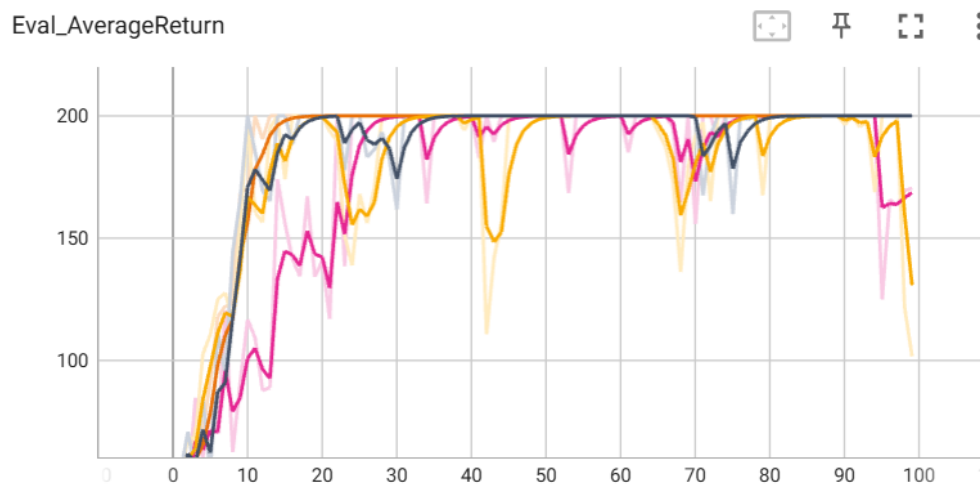


Figure 3: CartPole experiments with large batch size

Command line configurations:

```
# Small batch experiments
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
--exp_name cartpole
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
-rtg --exp_name cartpole_rtg
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
-na --exp_name cartpole_na
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
-rtg -na --exp_name cartpole_rtg_na
# Large batch experiments
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 4000 \
--exp_name cartpole_lb
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 4000 \
```

```
-rtg --exp_name cartpole_lb_rtg
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 4000 \
-na --exp_name cartpole_lb_na
python cs285/scripts/run_hw2.py --env_name CartPole-v0 -n 100 -b 4000 \
-rtg -na --exp_name cartpole_lb_rtg_na
```

According to the graphs, the reward-to-go estimator generally performs better than the trajectory-centric estimator, especially in the large batch experiments. Advantage normalization appears to help stabilize training, as seen in the smoother learning curves. The larger batch size leads to more stable updates and better performance, particularly in the large batch experiments.

The normalization of advantages helps to reduce variance in the policy gradient updates, leading to more stable learning.

The larger batch size does make a significant impact, as seen in the large batch experiments where the learning curves are smoother and the final performance is generally better compared to the small batch experiments.

5 Neural Network Baseline

- Plot a learning curve for the baseline loss.
- Plot a learning curve for the eval return. You should expect to achieve an average return over 300 for the baselined version.
- Run another experiment with a decreased number of baseline gradient steps (`-bgs`) and/or baseline learning rate (`-blr`). How does this affect (a) the baseline learning curve and (b) the performance of the policy?
- **Optional:** Add `-na` back to see how much it improves things. Also, set `video_log_freq 10`, then open TensorBoard and go to the “Images” tab to see some videos of your HalfCheetah walking along!

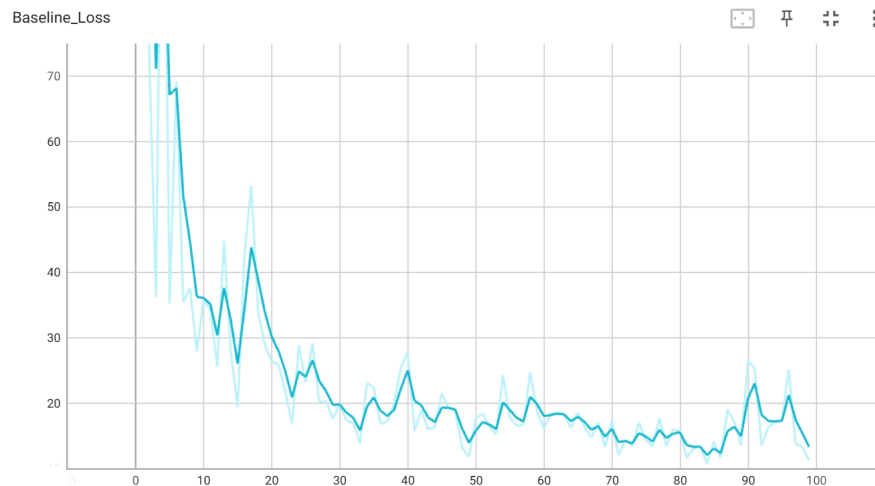


Figure 4: Learning curve for the baseline loss

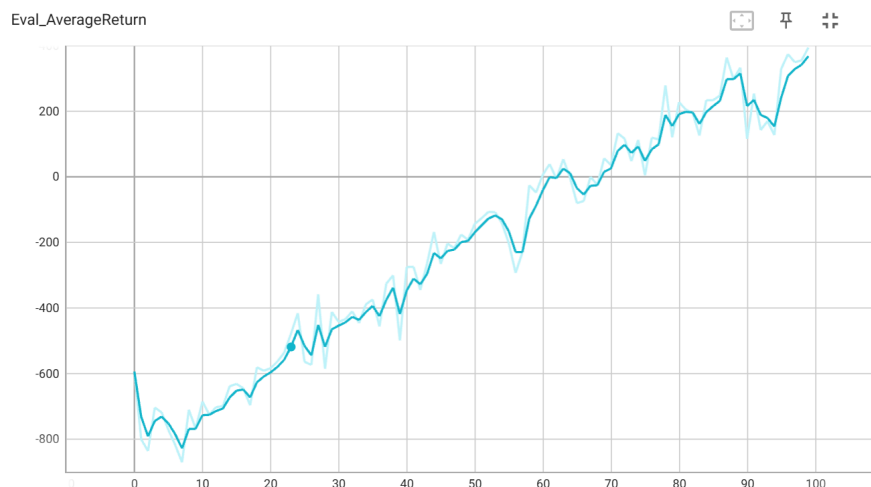


Figure 5: Learning curve for the baseline return

In the above experiments, we used the following default command line configurations:

```
python cs285/scripts/run_hw2.py --env_name HalfCheetah-v4 \
-n 100 -b 5000 -rtg --discount 0.95 -lr 0.01 \
```

```
--use_baseline -blr 0.01 -bgs 5 --exp_name cheetah_baseline
```

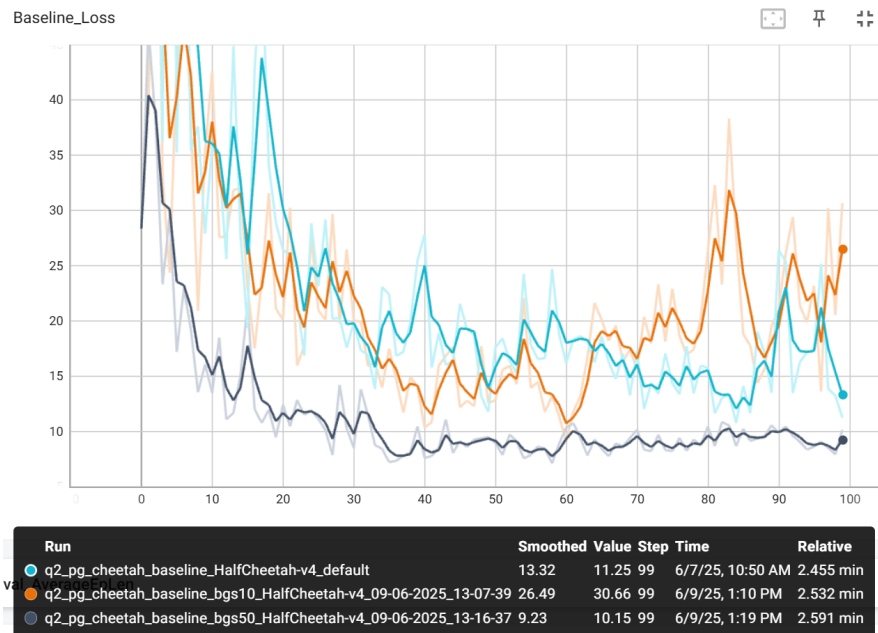


Figure 6: Baseline loss in different baseline gradient steps

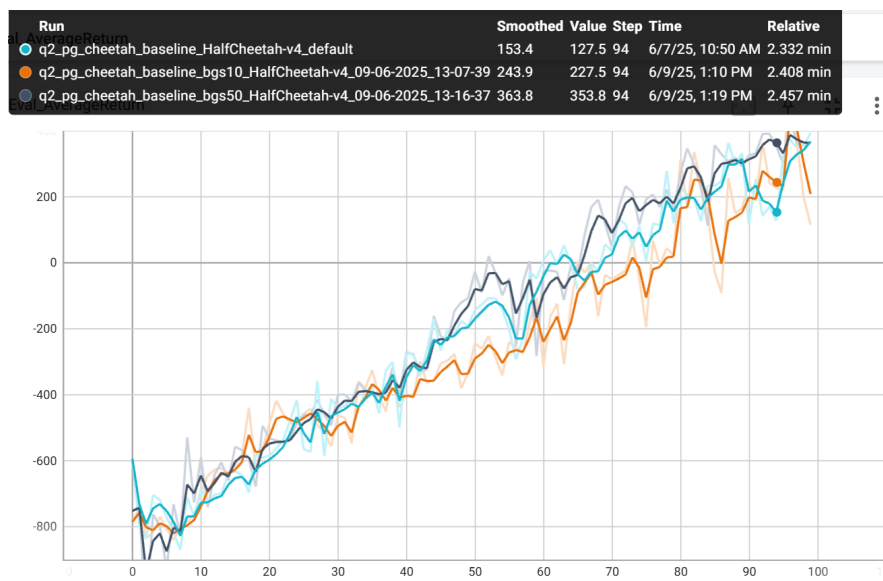


Figure 7: Baseline return in different baseline gradient steps

We try different baseline gradient steps 5, 10, and 50. The results show that increasing the number of baseline gradient steps leads to a more stable and smaller curve for the baseline loss, and the performance of the policy improves as well, achieving a higher average return.

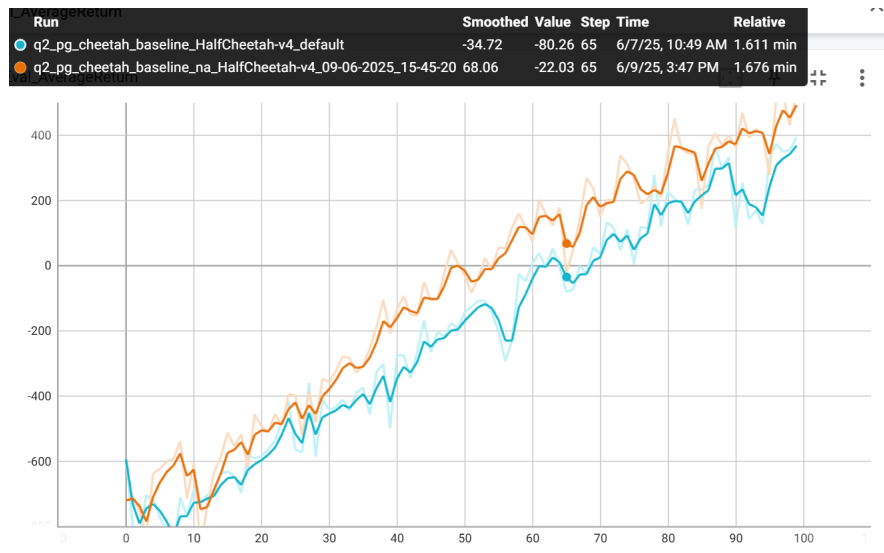


Figure 8: Learning curve for the baseline return

Compared to the default settings, adding `-na` (normalization) leads to a more stable learning curve and better performance of the policy, achieving a higher average return.

6 Generalized Advantage Estimation

- Provide a single plot with the learning curves for the **LunarLander-v2** experiments that you tried. Describe in words how λ affected task performance. The run with the best performance should achieve an average score close to 200 (180+).
- Consider the parameter λ . What does $\lambda = 0$ correspond to? What about $\lambda = 1$? Relate this to the task performance in **LunarLander-v2** in one or two sentences.

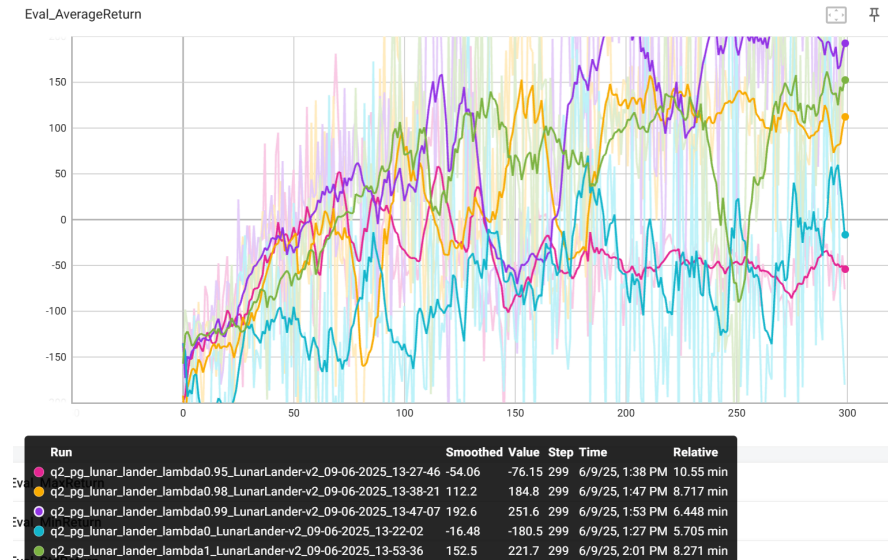


Figure 9: GAE Learning curve for the baseline return

$\lambda = 0$: Corresponds to using only the immediate rewards at each timestep (ignoring future rewards), effectively disabling Generalized Advantage Estimation (GAE). This results in lower bias but higher variance in advantage estimation, which may lead to unstable training.

$\lambda = 1$: Corresponds to fully incorporating future rewards (similar to Monte Carlo methods), using complete GAE. This reduces variance but introduces higher bias, potentially slowing down convergence.

7 Hyperparameter Tuning

1. Provide a set of hyperparameters that achieve high return on **InvertedPendulum-v4** in as few environment steps as possible.
2. Show learning curves for the average returns with your hyperparameters and with the default settings, with environment steps on the x -axis. Returns should be averaged over 5 seeds.

8 (Extra Credit) Humanoid

1. Plot a learning curve for the Humanoid-v4 environment. You should expect to achieve an average return of at least 600 by the end of training. Discuss what changes, if any, you made to complete this problem (for example: optimizations to the original code, hyperparameter changes, algorithmic changes).

9 Analysis

Consider the following infinite-horizon MDP:

$$a_1 \curvearrowright s_1 \xrightarrow{a_2} s_F$$

At each step, the agent stays in state s_1 and receives reward 1 if it takes action a_1 , and receives reward 0 and terminates the episode otherwise. Parametrize the policy as stationary (not dependent on time) with a single parameter:

$$\pi_\theta(a_1|s_1) = \theta, \pi_\theta(a_2|s_1) = 1 - \theta$$

1. Applying policy gradients

- (a) Use policy gradients to compute the gradient of the expected return $R(\tau)$ with respect to the parameter θ . **Do not use discounting.**

Hint: to compute $\sum_{k=1}^{\infty} k\alpha^{k-1}$, you can write:

$$\sum_{k=1}^{\infty} k\alpha^{k-1} = \sum_{k=1}^{\infty} \frac{d}{d\alpha} \alpha^k = \frac{d}{d\alpha} \sum_{k=1}^{\infty} \alpha^k$$

Solution:

- (b) Compute the expected return of the policy $\mathbb{E}_{\tau \sim \pi_\theta} R(\tau)$ directly. Compute the gradient of this expression with respect to θ and verify that this matches the policy gradient.

Solution:

2. Compute the variance of the policy gradient in closed form and describe the properties of the variance with respect to θ . For what value(s) of θ is variance minimal? Maximal? (Once you have an exact expression for the variance you can eyeball the min/max).

Hint: Once you have it expressed as a sum of terms $P(\theta)/Q(\theta)$ where P and Q are polynomials, you can use a symbolic computing program (Mathematica, SymPy, etc) to simplify to a single rational expression.

Solution:

3. Apply return-to-go as an advantage estimator.

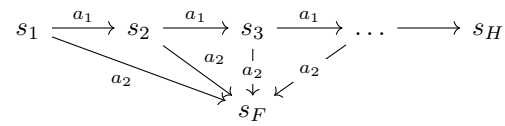
- (a) Write the modified policy gradient and confirm that it is unbiased.

Solution:

- (b) Compute the variance of the return-to-go policy gradient and plot it on $[0, 1]$ alongside the variance of the original estimator.

Solution:

4. Consider a finite-horizon H -step MDP with sparse reward:



The agent receives reward R_{\max} if it arrives at s_H and reward 0 if it arrives at s_F (a terminal state). In other words, the return for a trajectory τ is given by:

$$R(\tau) = \begin{cases} 1 & \tau \text{ ends at } s_H \\ 0 & \tau \text{ ends at } s_F \end{cases}$$

Using the same policy parametrization as above, consider off-policy policy gradients via importance sampling. Assume we want to compute policy gradients for a policy π_θ with samples drawn from $\pi_{\theta'}$.

- Write the policy gradient with importance sampling.
- Compute its variance.

10 Survey

Please estimate, in minutes, for each problem, how much time you spent (a) writing code and (b) waiting for the results. This will help us calibrate the difficulty for future homeworks.

- **Policy Gradients:**
- **Neural Network Baseline:**
- **Generalized Advantage Estimation:**
- **Hyperparameters and Sample Efficiency:**
- **Humanoid:**
- **Humanoid:**
- **Analysis – applying policy gradients:**
- **Analysis – PG variance:**
- **Analysis – return-to-go:**
- **Analysis – importance sampling:**