

# 1 CartPole

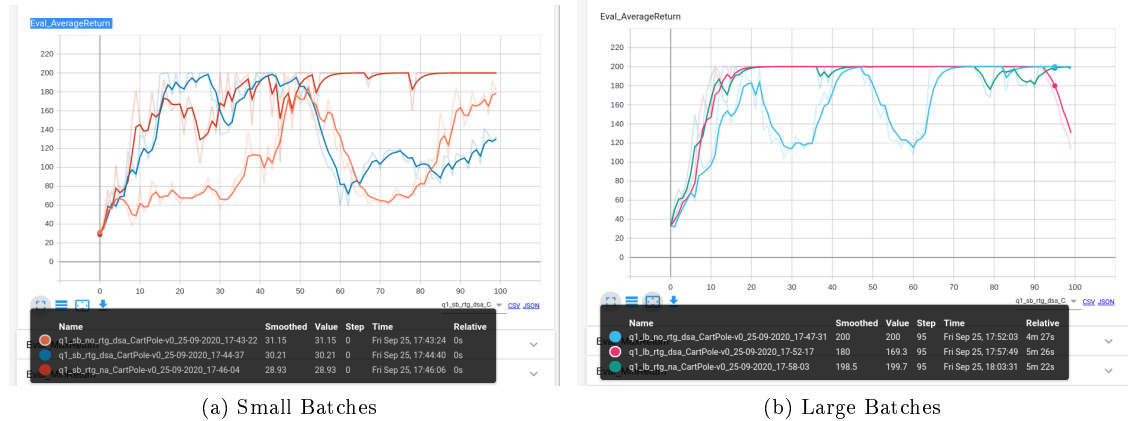


Figure 1: CartPole Experiment

**Which value estimator has better performance without advantage-standardization: the trajectorycentric one, or the one using reward-to-go?**

The plots with names containing “dsa” have advantage normalization disabled. In figure (a) this included the orange and blue plots. The blue plot which represent the “reward to go” trajectory, trains much more quickly initially but then fails to converge and drops off later. The orange plot trains more slowly and also stays at the peak for less time. Although it seems to get a better recovery the second time around, I think the performance of the blue plot which uses reward to go is better. When we look at the larger batches in figure (b) it’s super clear that the reward to go dominates the trajectory centric learner. Thus I conclude that without advantage normalization reward to go has a performance that is confidently superior.

**Did advantage standardization help?**

The beneficial effect of normalization is most visible with small batches. The red plot is clearly the best performer. In figure (b) the green plot which represents the standardized advantage performs nearly as good as the red plot. I think that they’re close enough to call a tie.

### Did the batch size make an impact?

An increase in batch size led to performance enhancements across the board. I would argue it was the most important parameter.

**Provide the exact command line configurations (or #@params settings in Colab) you used to run your experiments, including any parameters changed from their defaults.**

```
Experiment 1 (CartPole)

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -dsa --exp_name q1_sb_no_rtg_dsa

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -rtg -dsa --exp_name q1_sb_rtg_dsa

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 1000 \
    -rtg --exp_name q1_sb_rtg_na

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -dsa --exp_name q1_lb_no_rtg_dsa

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -rtg -dsa --exp_name q1_lb_rtg_dsa

[ ]: !python ./run_hw2.py --env_name CartPole-v0 -n 100 -b 5000 \
    -rtg --exp_name q1_lb_rtg_na
```

Figure 2: CartPole section of cs285/scripts/run\_experiments.ipynb

## 2 Inverted Pendulum

Given the  $b^*$  and  $r^*$  you found, provide a learning curve where the policy gets to optimum (maximum score of 1000) in less than 100 iterations.

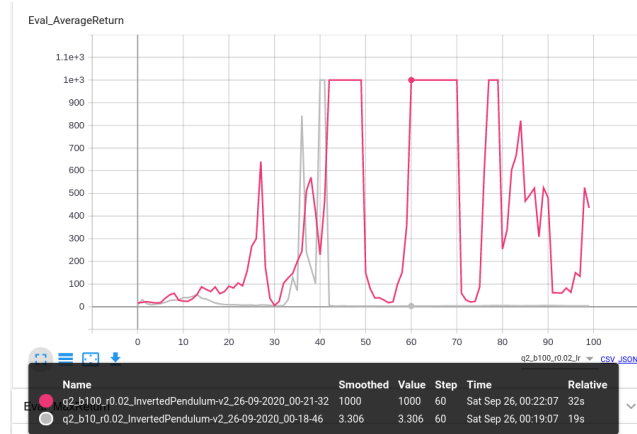


Figure 3: Pendulum Experiment

I included two runs because although using a batch size of 10 and a learning rate of 0.02 (the grey line) technically reaches 1000, it barely manages and then it crashes and never comes close again. The pink line on the other hand ( $b=100$ ,  $lr=0.02$ ) reached 1000 mark pretty consistently relative to how much fluctuation there was for every run in the pendulum experiment.

Provide the exact command line configurations you used to run your experiments

```

Experiment 2 (Inverted Pendulum)

[ ]: def run_pendulum(r, b):
    !python ./run_hw2.py --env_name InvertedPendulum-v2 \
    --ep_len 1000 --discount 0.9 -n 100 -l 2 -s 64 -b {b} -lr {r} -rtg \
    --exp_name q2_b{b}_r{r}

    for b in [10, 100, 200, 500, 800, 1000, 1500, 3000]:
        for lr in [0.005, .01, .02, .05, .1]:
            run_pendulum(lr, b)

```

Figure 4: Pendulum section of cs285/scripts/run\_experiments.ipynb

### 3 Lunar Landing

Plot a learning curve for the above command. You should expect to achieve an average return of around 180 by the end of training.

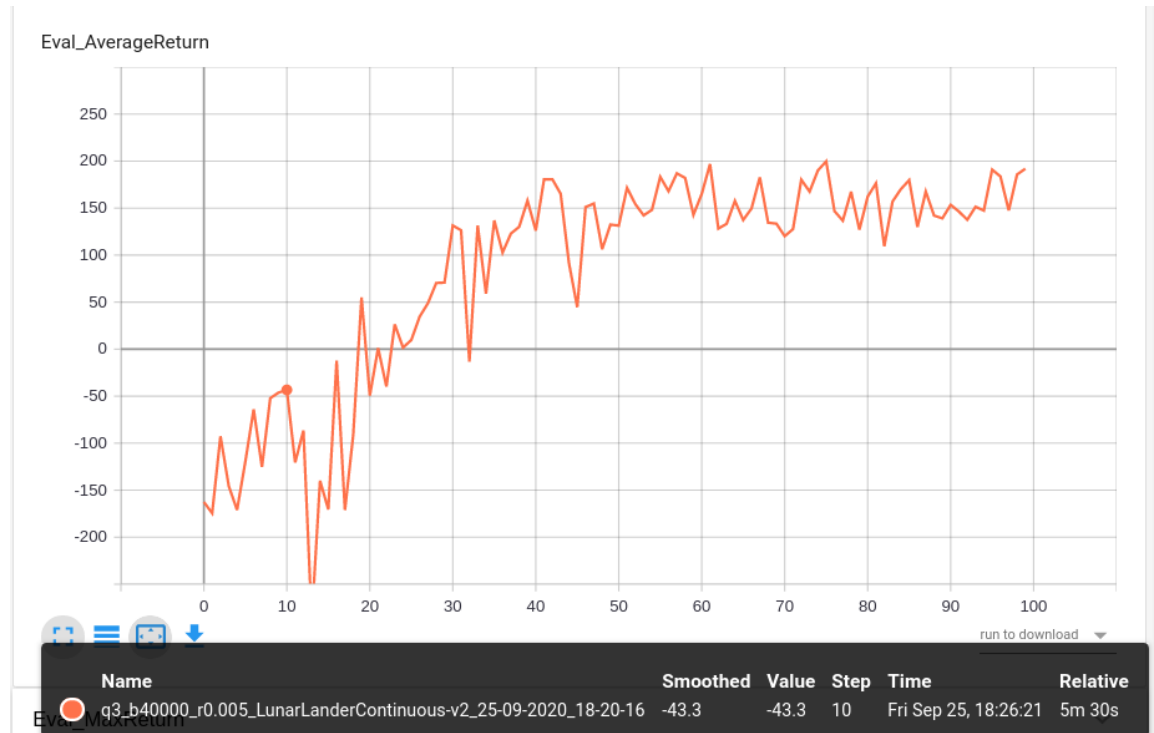


Figure 5: Lunar Landing Experiment

```
Experiment 3 (LunarLander)

!python ./run_hw2.py --env_name LunarLanderContinuous-v2 --ep_len 1000
--discount 0.99 -n 100 -l 2 -s 64 -b 40000 -lr 0.005 --reward_to_go
--nn_baseline --exp_name q3_b40000_r0.005
```

Figure 6: Lunar Landing section of cs285/scripts/run\_experiments.ipynb

## 4 HalfCheetah

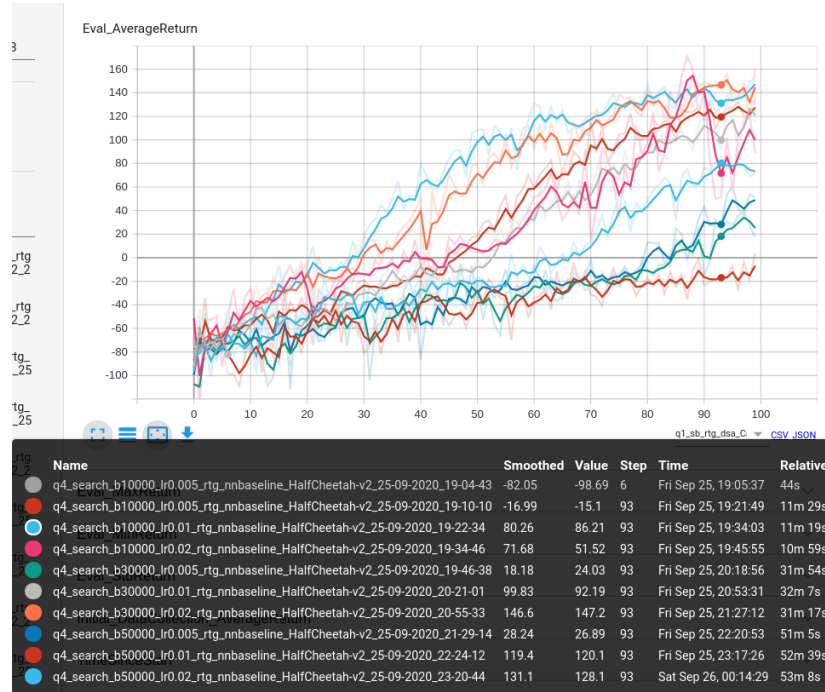


Figure 7: Cheetah finding batch size and learning rate

**Provide a single plot with the learning curves for the HalfCheetah experiments that you tried. Describe in words how the batch size and learning rate affected task performance.**

The learning rate of 0.02 was the most effective. At one point in the graph, the 3 highest values all have learning rate of 0.02. The lowest runs all had the smallest learning rate so there is definitely some correlation between learning rate and progress. Seems like the larger the learning rate the better the learner did. The same goes for batch sizes, the larger the better.

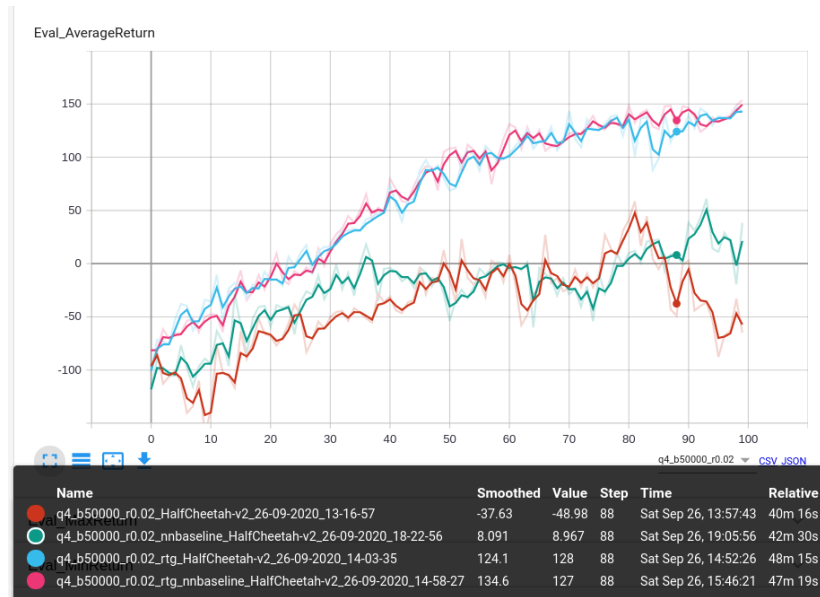


Figure 8: Cheeta Additional Experiments

### Experiment 4 (HalfCheeta)

```
[ ]: def run_cheeta(r, b):
    !python ./run_hw2.py --env_name HalfCheetah-v2 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b {b} -lr {r} -rtg --nn_baseline \
    --exp_name q4_search_b{b}_lr{r}_rtg_nnbaseline

    for b in [10000, 30000, 50000]:
        for lr in [0.005, 0.01, 0.02]:
            run_lunar(lr, b)
```

Now that we have selected b=50000, lr=0.02... Running additional experiments:

```
[ ]: !python ./run_hw2.py --env_name HalfCheetah-v2 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b 50000 -lr 0.02 \
    --exp_name q4_b50000_r0.02

[ ]: !python ./run_hw2.py --env_name HalfCheetah-v2 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b 50000 -lr 0.02 -rtg \
    --exp_name q4_b50000_r0.02_rtg

[ ]: !python ./run_hw2.py --env_name HalfCheetah-v2 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b 50000 -lr 0.02 --nn_baseline \
    --exp_name q4_b50000_r0.02_nnbaseline

[ ]: !python ./run_hw2.py --env_name HalfCheetah-v2 --ep_len 150 \
    --discount 0.95 -n 100 -l 2 -s 32 -b 50000 -lr 0.02 -rtg --nn_baseline \
    --exp_name q4_b50000_r0.02_rtg_nnbaseline
```

Figure 9: Cheeta section of cs285/scripts/run\_experiments.ipynb