

### 1.1 TD-Learning Bias (2 points)

We say an estimator  $f_D$  of  $f$  constructed using data  $\mathcal{D}$  sampled from process  $P$  is *unbiased* when  $\mathbb{E}_{\mathcal{D} \sim P}[f_D(x)] - f(x) = 0$  at each  $x$ .

Assume  $\hat{Q}$  is a noisy (but unbiased) estimate for  $Q$ . Is the Bellman backup  $B\hat{Q} = r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')$  an unbiased estimate of  $BQ$ ?

☐ Yes ☒ No **NO**  $E[r(s, a)] = r(s, a)$  so  $r(s, a)$  doesn't introduce bias  
check if  $E[B\hat{Q}] = BQ$ .  $E[\max \hat{Q}] \neq \max Q$  since  $\hat{Q}$  is a noisy estimate

	I.	II.	III.
1. $N = 1$ and ...	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
(a) on-policy in tabular setting	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
(b) off-policy in tabular setting	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
2. $N > 1$ and ...	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
(a) on-policy in tabular setting	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
(b) off-policy in tabular setting	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
3. In the limit as $N \rightarrow \infty$ (no bootstrapping) ...	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
(a) on-policy in tabular setting	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
(b) off-policy in tabular setting	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

### 1.3 Variance of $Q$ Estimate (2 points)

Which of the three cases ( $N = 1$ ,  $N > 1$ ,  $N \rightarrow \infty$ ) would you expect to have the highest-variance estimate of  $Q$  for fixed dataset size  $B$  in the limit of infinite iterations  $k$ ? Lowest-variance?

Highest variance: ☐  $N = 1$  ☒  $N > 1$  ☒  $N \rightarrow \infty$   
Lowest variance: ☒  $N = 1$  ☐  $N > 1$  ☐  $N \rightarrow \infty$

Figure 1: Enter Caption

Figure 1: Enter Caption

### 1.4 Function Approximation (2 points)

Now say we want to represent  $Q$  via function approximation rather than with a tabular representation. Assume that for any deterministic policy  $\pi$  (including the optimal policy  $\pi^*$ ), function approximation can represent the true  $Q^\pi$  exactly. Which of the following statements are true?

- ☐ When  $N = 1$ ,  $Q_{\phi_{k+1}}$  is an unbiased estimate of the  $Q$ -function of the last policy  $Q^{\pi_k}$ .
- ☐ When  $N = 1$  and in the limit as  $B \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $Q_{\phi_k}$  converges to  $Q^*$ .
- ☒ When  $N > 1$  (but finite) and in the limit as  $B \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $Q_{\phi_k}$  converges to  $Q^*$ .
- ☒ When  $N \rightarrow \infty$  and in the limit as  $B \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $Q_{\phi_k}$  converges to  $Q^*$ .

Figure 2: Enter Caption

Figure 2: Enter Caption

$$\begin{aligned}
\phi_{k+1} &\leftarrow \underset{\phi \in \Phi}{\operatorname{argmin}} \sum_{j,t} (y_{j,t} - Q_{\phi}(s_{j,t}, a_{j,t})) \\
y_{j,t} &= \sum_{t'=t}^{t+N-1} \gamma^{t'-t} r_{j,t'} + \gamma^N \max_{a_{j,t+N}} \sum_{a_{j,t+N}}^T E_{\pi} [r(s_{j,t}, a_{j,t}) | s_{j,t}, a_{j,t}] \\
&= \sum_{t'=t}^{t+N-1} \gamma^{t'-t} r_{j,t'} + \gamma^N \max_{a_{j,t+N}} \sum_{a_{j,t+N}}^T E_{\pi} \left[ \frac{\prod_{t'=t}^{t+N-1} \pi(a_{j,t'+1})}{\prod_{t'=t}^T \pi(a_{j,t'+1})} r(s_{j,t}, a_{j,t}) \right] \\
&= \sum_{t'=t}^{t+N-1} \gamma^{t'-t} r_{j,t'} + \frac{P_{\phi}(\tau)}{P_{\pi}(\tau)} \cdot \gamma^N \max_{a_{j,t+N}} Q_{\pi}(s_{j,t+N}, a_{j,t+N}) \\
Q_{\phi}(s_{j,t}, a_{j,t}) &= \sum_{j,t} E_{\pi_{\phi}} [r(s_{j,t}, a_{j,t}) | s_{j,t}, a_{j,t}] \\
&= \sum_{j,t} \frac{P_{\phi}(a_{j,t} | s_{j,t})}{P_{\pi}(a_{j,t} | s_{j,t})} \cdot E_{\pi} [r(s_{j,t}, a_{j,t}) | s_{j,t}, a_{j,t}] \\
&= \frac{P_{\phi}(\tau)}{P_{\pi}(\tau)} \cdot Q_{\pi}(s_{j,t}, a_{j,t})
\end{aligned}$$

Figure 3: 1.5

Figure 3: Enter Caption

What if  $N=1$  or  $N \rightarrow \infty$

a.  $N=1$  Case, all  $Q$  learning

plug into  $\phi_{k+1} \leftarrow$

$$\begin{aligned}
y_{j,t} &\leftarrow r_{j,t} + \gamma \max_{a_{j,t+1}} \sum_{t'=t+1}^T E_{\pi} \left[ \frac{P_{\phi}(a_{j,t+1} | s_{j,t})}{P_{\pi}(a_{j,t+1} | s_{j,t})} r(s_{j,t+1}, a_{j,t+1}) \right] \\
&\leftarrow r_{j,t} + \gamma \max_{a_{j,t+1}} Q_{\phi}(s_{j,t+1}, a_{j,t+1}) \cdot \frac{P_{\phi}(\tau)}{P_{\pi}(\tau)} \\
Q_{\phi}(s_{j,t}, a_{j,t}) &= Q_{\phi}(s_{j,t}, a_{j,t}) \cdot \frac{P_{\phi}(\tau)}{P_{\pi}(\tau)}
\end{aligned}$$

b.  $N = \infty$  case

$$\phi_{k+1} \leftarrow \underset{\phi \in \Phi}{\operatorname{argmin}} \sum_{j,t} \left( \sum_{t'=t}^{\infty} \gamma^{t'-t} r_{j,t'} - \frac{P_{\phi}(\tau)}{P_{\pi}(\tau)} Q_{\pi}(s_{j,t}, a_{j,t}) \right)^2$$

Figure 4: Enter Caption

## 1 2.4

```

rm -rf ./data/*
python cs285/scripts/run_hw3_dqn.py -cfg experiments/dqn/cartpole.yaml
python cs285/scripts/graph_results.py

```

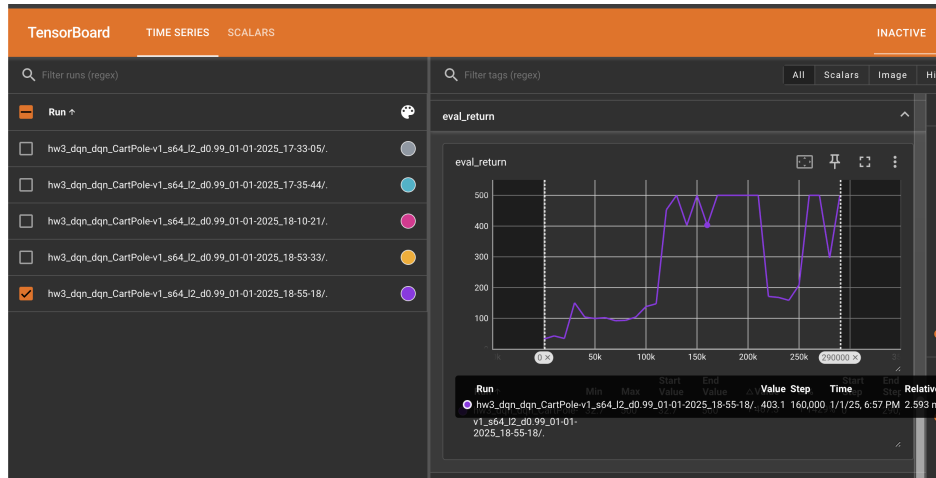


Figure 5: Enter Caption

**Note:** Purple is seed 3, blue is seed 2, grey is seed 1.

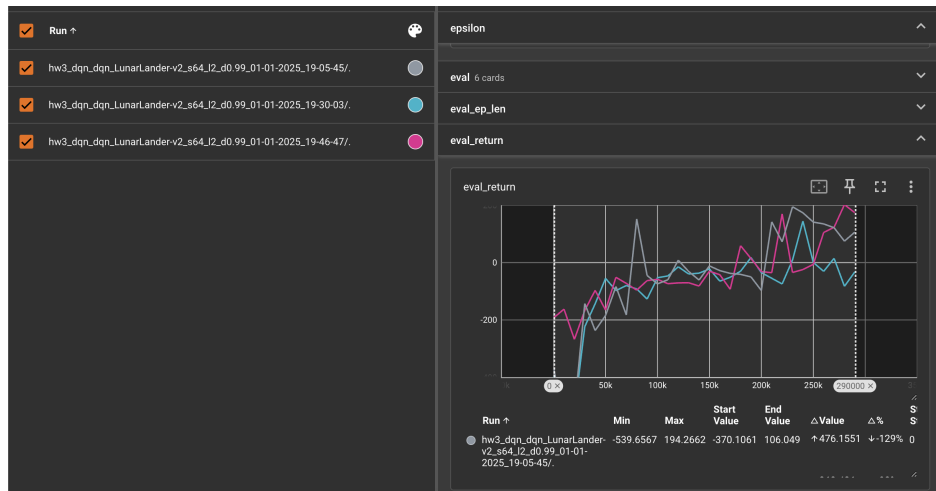


Figure 6: Enter Caption

Change the learning rate to 0.05 in the YAML config file. Analyze the impact on:

1. **Predicted Q-values:** Predicted Q-values increased but became more volatile.
2. **Critic error:** Critic error increased.

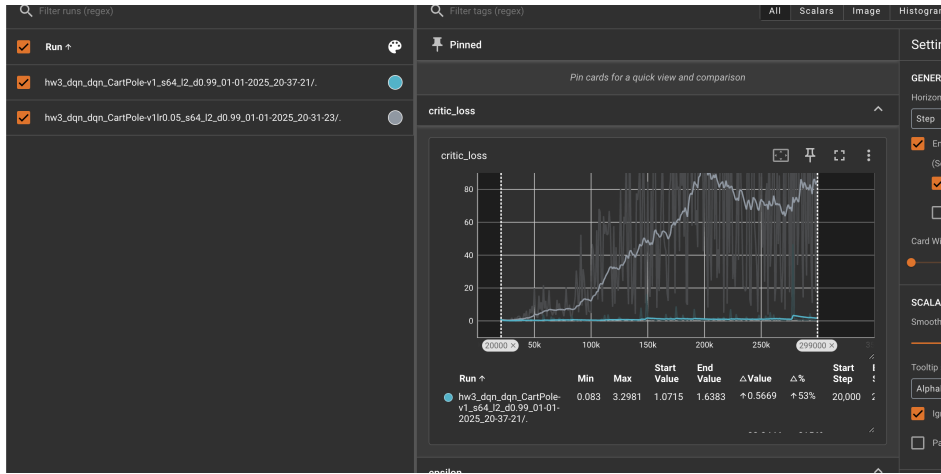


Figure 7: Increased learning rate and critic loss

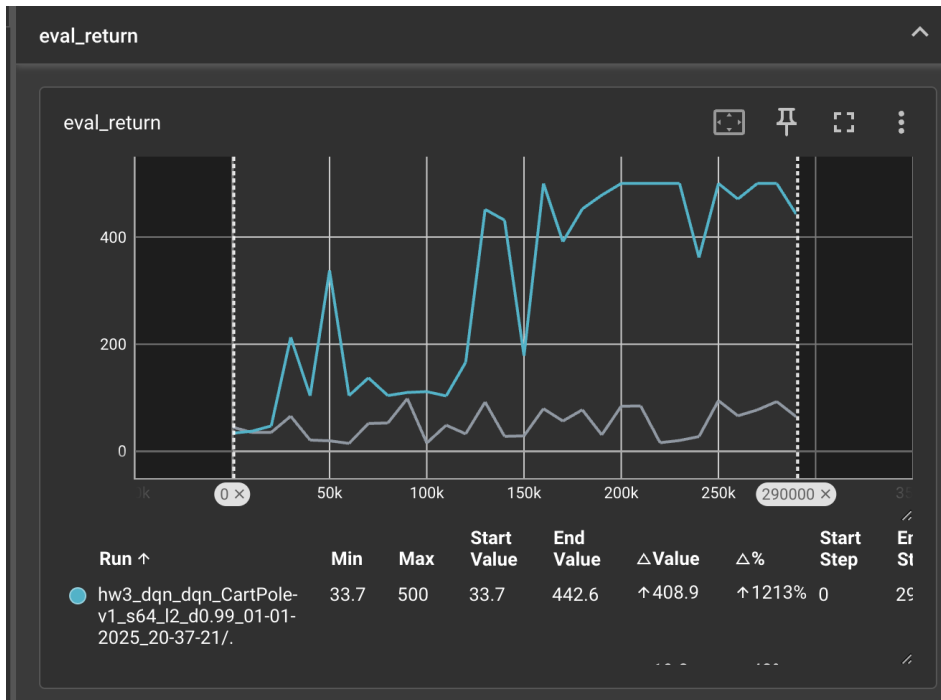


Figure 8:  $lr = 0.05$  and eval return

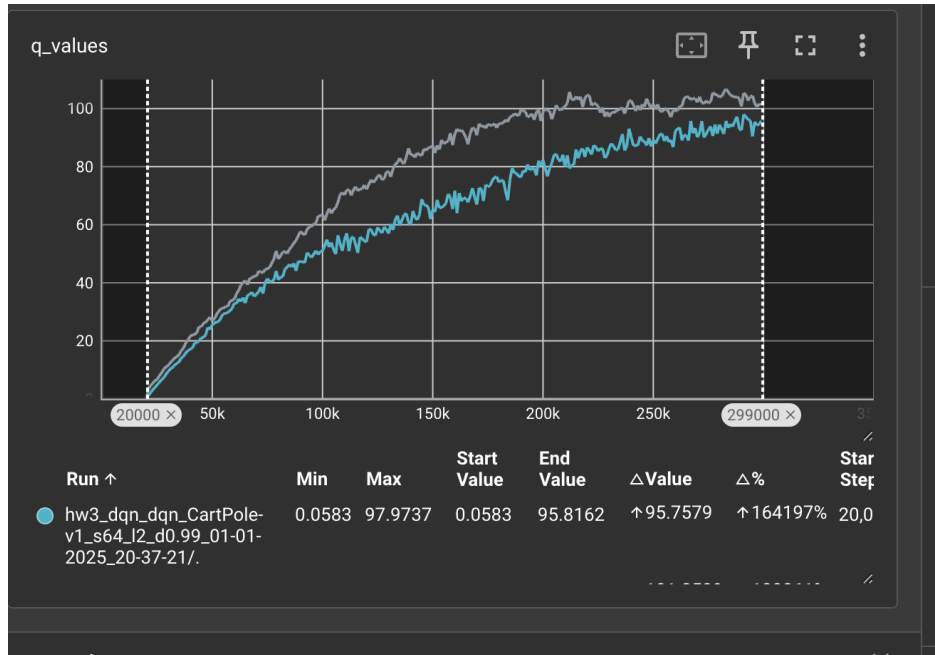


Figure 9:  $lr = 0.05$  and q values

**Explanation:** The predicted Q-values increased because the higher learning rate led to more aggressive updates via the gradient function, trending toward favoring higher future rewards. However, the increased volatility caused higher variance in the predicted Q-values.

Since the Q-values became noisier, the Bellman error also increased. This is because the Bellman error depends on the estimation of future rewards using the Q-function for the target.

Higher variance in Q-values contributes to higher variance in the Bellman error. A high learning rate causes the neural network to overshoot the optimal values of the target critic Q-values, leading to increased critic error.

## 2 2.5



Figure 10: DQN and eval return

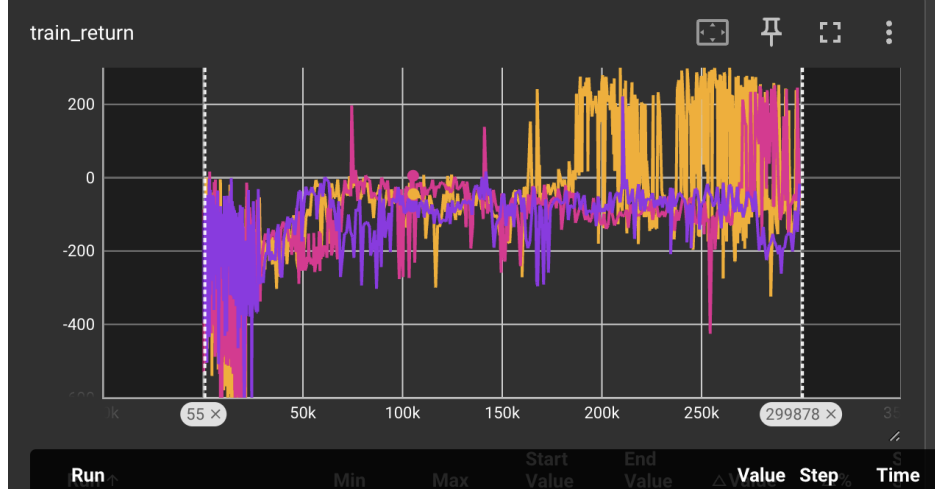


Figure 11:  $lr = 0.05$  and  $q$  values

Here the double  $q$  learning seems to take longer to reach eval return around 200 compared to the vanilla approach since it could be possible that our two critics take longer to be trained to predict  $Q$  values compared to only needing to train one critic for both choosing actions and evaluating the  $q$  value for the selected action.

### 3 2.6

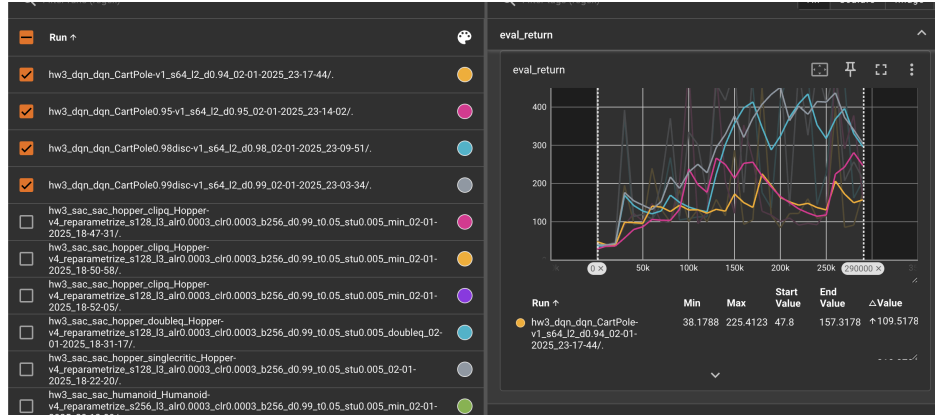


Figure 12: Discount factor vs eval return for cartpole

## Impact of Discount Factor

The discount factor determines the extent to which the DQN agent prioritizes future rewards over present rewards, i.e., it affects the  $Q$ -value. A higher discount factor results in a higher evaluation return, as the DQN agent assigns more importance to future rewards, weighing them almost equally with current rewards.

By weighing future rewards higher,  $Q$ -value estimations align more closely with the target values, improving performance in tasks where long-term outcomes are critical. This leads to better policy learning that balances immediate and delayed rewards effectively.

### 4 3.1.1

```
python cs285/scripts/run_hw3_sac.py -cfg experiments/sac/sanity_pendulum.yaml
```



Figure 13:  $lr = 0.05$  and  $q$  values

### 5 3.1.3



Figure 14:  $lr = 0.05$  and  $q$  values

Return is significantly higher when using Reinforce10 (yellow graph). Using only one sample to train the gradient introduces bias toward the training sample, causing the Q-values to deviate further from their expected values.

### 6 3.1.4

```
python cs285/scripts/run_hw3_sac.py -cfg experiments/sac/halfcheetah_reparametrize.yaml
python cs285/scripts/run_hw3_sac.py -cfg experiments/sac/sanity_invertedpendulum_reinforce.yaml
python cs285/scripts/run_hw3_sac.py -cfg experiments/sac/halfcheetah_reinforce1.yaml
```



Figure 15:  $lr = 0.05$  and  $q$  values