

Exercise 10

Kai Schneider

July 6, 2021

Task 1 DQN on the Cart-Pole

a.)

- NNs tend to diverge (due to high correlation states and actions, NOT i.i.d)
→ Experience Replay (sampling training data from buffer of past experiences)
- oscillation due to big changes in Q-values
→ much lower learning rate than in other DL-scenarios
→ reward clipping, normalize $[0, 1]$
- bias introduced at beginning of training and instability due to similar subsequent steps
→ 2 NNs. Predictions of a target NN are used to get the Q-values during training. This target network is only synced periodically. The main NN then uses this to backpropagate its weights.