

RL: Deep DQN

Marius Lindauer

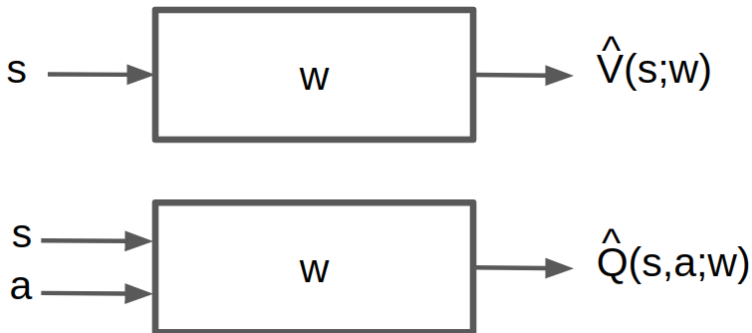


Winter Term 2021

RL with Function Approximation

- Represent state-action value function by Q -network with weights \vec{w}

$$\hat{Q}(s, a; \vec{w}) \approx Q(s, a)$$



Recall: Incremental Model-Free Control Approaches

- ▶ Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value
- ▶ In Monte Carlo methods, use a return G_t as a substitute target

$$\Delta \vec{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s_t, a_t; \vec{w})$$

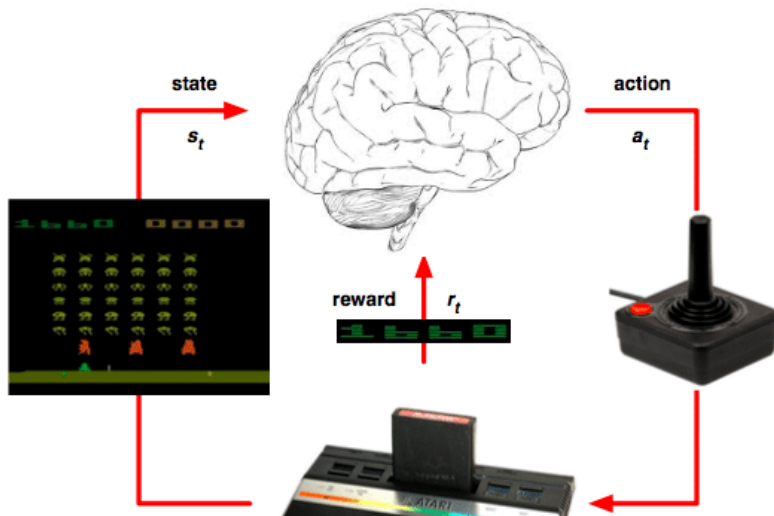
- ▶ For SARSA instead use a TD target $r + \gamma \hat{Q}(s', a'; \vec{w})$ which leverages the current function approximations value

$$\Delta \vec{w} = \alpha(r + \gamma \hat{Q}(s', a'; \vec{w}) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

- ▶ For Q-learning instead use a TD target $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w})$ which leverages the max of the current function approximations value

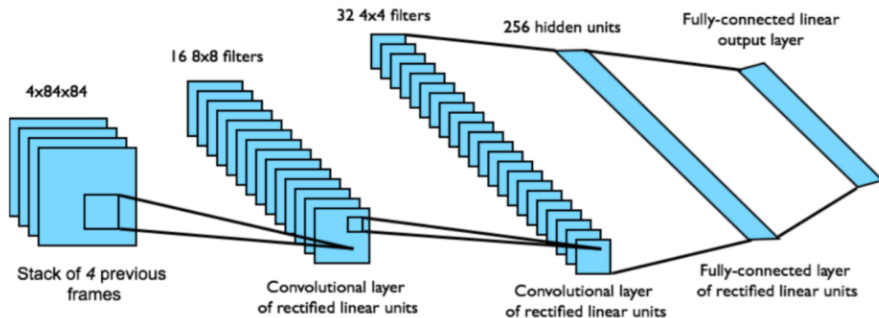
$$\Delta \vec{w} = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

Using these Ideas to do Deep RL in Atari



Using these Ideas to do Deep RL in Atari

- ▶ End-to-end learning of values $Q(s, a)$ from pixels s
- ▶ Input state s is stack of raw pixels from last 4 frames
- ▶ Output is $Q(s, a)$ for 18 joystick/button positions
- ▶ Reward is change in score for that step
- ▶ Network architecture and hyperparameters fixed across all games



Q-Learning with Value Function Approximation

- ▶ Minimize MSE loss by stochastic gradient descent
- ▶ Converges to the optimal $Q^*(s, a)$ using **table lookup** representation
- ▶ But Q-learning with VFA can diverge
- ▶ Two of the issues causing problems:
 - ▶ Correlations between samples violates i.i.d assumption of DNNs
 - ▶ Non-stationary targets
- ▶ Deep Q-learning (DQN) addresses both of these challenges by
 - ▶ Experience replay
 - ▶ Fixed Q-targets

DQNs: Replay Buffer

- ▶ To help remove correlations, store dataset (called a **replay buffer**) \mathcal{D} from prior experience
- ▶ To perform experience replay, repeat the following:
 1. $(s, a, r, s') \sim \mathcal{D}$: sample experience tuple from the dataset
 2. Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w})$
 3. Use stochastic gradient descent to update the network weights

$$\Delta \vec{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

DQNs: Replay Buffer

- ▶ To help remove correlations, store dataset (called a **replay buffer**) \mathcal{D} from prior experience
- ▶ To perform experience replay, repeat the following:
 1. $(s, a, r, s') \sim \mathcal{D}$: sample experience tuple from the dataset
 2. Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w})$
 3. Use stochastic gradient descent to update the network weights

$$\Delta \vec{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

- ▶ Remarks:
 - ▶ Fixed sized buffer \leadsto first-in-first-out scheme (as default implementation)
 - ▶ heuristic trade-off between performing new episodes and sampling from the replay buffer

DQNs: Replay Buffer

- ▶ To help remove correlations, store dataset (called a **replay buffer**) \mathcal{D} from prior experience
- ▶ To perform experience replay, repeat the following:
 1. $(s, a, r, s') \sim \mathcal{D}$: sample experience tuple from the dataset
 2. Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w})$
 3. Use stochastic gradient descent to update the network weights

$$\Delta \vec{w} = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

- ▶ Remarks:
 - ▶ Fixed sized buffer \leadsto first-in–first-out scheme (as default implementation)
 - ▶ heuristic trade-off between performing new episodes and sampling from the replay buffer
- \leadsto Can treat the target as a scalar, but the weights will get updated on the next round, changing the target value

DQNs: Fixed Q-Targets

- ▶ To help improve stability, fix the **target weights** used in the target calculation for multiple updates
- ▶ Target network uses a different set of weights than the weights being updated
- ▶ Let parameters \vec{w}^- be the set of weights used in the target and \vec{w} be the weights that are being updated
- ▶ Slight change to computation of target value:
 - ▶ $(s, a, r, s') \sim \mathcal{D}$: sample experience tuple from the dataset
 - ▶ Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}^-)$
 - ▶ Use stochastic gradient descent to update the network weights

$$\Delta \vec{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}^-) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

DQNs: Fixed Q-Targets

- ▶ To help improve stability, fix the **target weights** used in the target calculation for multiple updates
- ▶ Target network uses a different set of weights than the weights being updated
- ▶ Let parameters \vec{w}^- be the set of weights used in the target and \vec{w} be the weights that are being updated
- ▶ Slight change to computation of target value:
 - ▶ $(s, a, r, s') \sim \mathcal{D}$: sample experience tuple from the dataset
 - ▶ Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}^-)$
 - ▶ Use stochastic gradient descent to update the network weights

$$\Delta \vec{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \vec{w}^-) - \hat{Q}(s, a; \vec{w})) \nabla_{\vec{w}} \hat{Q}(s, a; \vec{w})$$

- ▶ Remark:
 - ▶ Hyperparameter how often you update \vec{w}^-
 - ▶ Trade-off between updating too often (\leadsto instability) and too rarely (\leadsto too old state information)

DQN Summary

- ▶ DQN uses experience replay and fixed Q-targets
- ▶ Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory \mathcal{D}
- ▶ Sample random mini-batch of transitions (s, a, r, s') from \mathcal{D}
- ▶ Compute Q-learning targets wrt old, fixed parameters \vec{w}^-
 - ▶ Update \vec{w}^- from time to time
- ▶ Optimizes MSE between Q-network and Q-learning targets
- ▶ Uses stochastic gradient descent