

Deep Q learning

Review

Model - Based

ML MB RL

Dyna

Prioritized Sweeping

Bayesian RL

Model-Free

Sarsa

Q-learning

On-Policy

Off-Policy

Eligibility Traces
Double Q-Learning

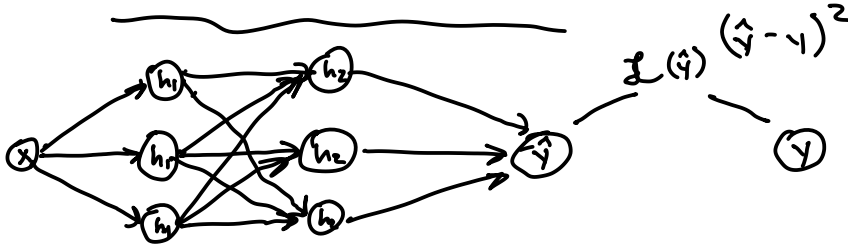
Classical/
Traditional
RL

1. Ex/Ex
2. Credit Assignment
3. Generalization ←

Function Approx

$$Q_{\theta}(s, a)$$

Neural Networks



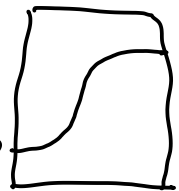
$$i_1 = w_1^T x + b_1$$

$$h_1 = \sigma(w_1^T x + b_1)$$

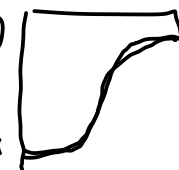
$$\hat{y} = w_3^T \sigma_2(w_2^T \sigma_1(w_1^T x + b_1) + b_2) + b_3$$

$$\hat{y} = W_3 \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) + b_3$$

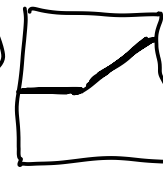
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



$$\sigma(x) = \tanh(x)$$



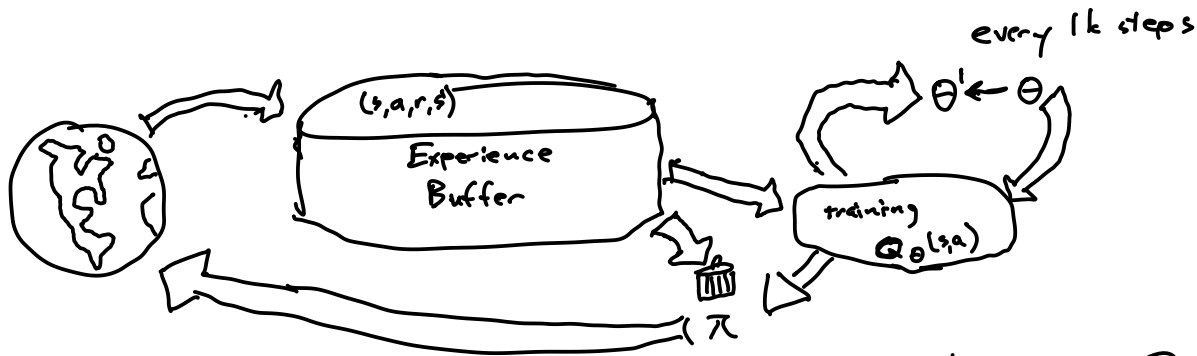
$$\sigma(x) = \max(0, x)$$



1. Soln. to 1.: Freeze target $\theta \rightarrow \theta'$

$$y = r + \gamma \max_{a'} Q_{\theta'}(s', a')$$

Soln. to 2+3: Use experience replay buffer



DQN

$$(r + \gamma \max_{a'} Q_{\theta'}(s, a) - Q_{\theta}(s, a))^2$$

- Double DQN

$$(r + \gamma Q_{\theta'}(s', \arg\max_{a'} Q_{\theta}(s', a')) - Q_{\theta}(s, a))^2$$

= Prioritized Replay

(s, a, r, s') chosen in proportion to size of last TD error

- Dueling Networks

← value

$$Q(s, a) = A(s, a) + V(s)$$

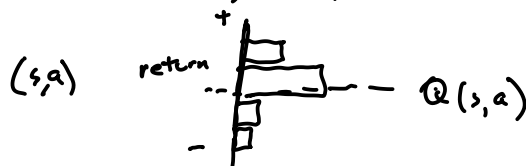
← advantage

- n-step Multistep learning

$$(r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n \max_{a'} Q_{\theta'}(s_{t+n}, a') - Q_{\theta}(s_t, a_t))^2$$

- Distributional RL

instead of learning Q , learn distribution of returns



- Noisy Nets

add noise in Q network

Rainbow