

Deep Reinforcement Learning

[Human-Level Control through deep reinforcement learning, Nature 2015]

CS 486/686

University of Waterloo

Lecture 20: July 10, 2017

Outline

- Value Function Approximation
 - Linear approximation
 - Neural network approximation
 - Deep Q-network

Quick recap

- Markov Decision Processes: value iteration

$$V(s) \leftarrow \max_a R(s) + \gamma \sum_{s'} \Pr(s'|s, a) V(s')$$

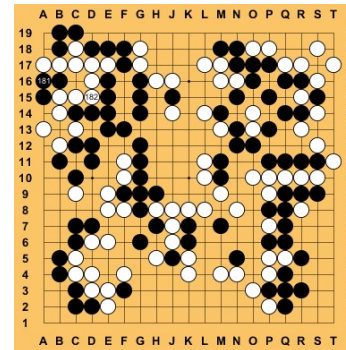
- Reinforcement Learning: Q-Learning

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

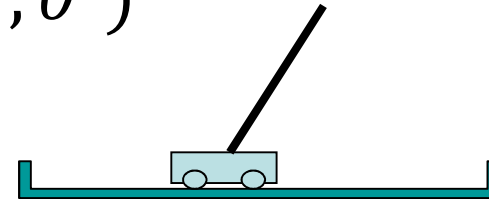
- Complexity depends on number of states and actions

Large State Spaces

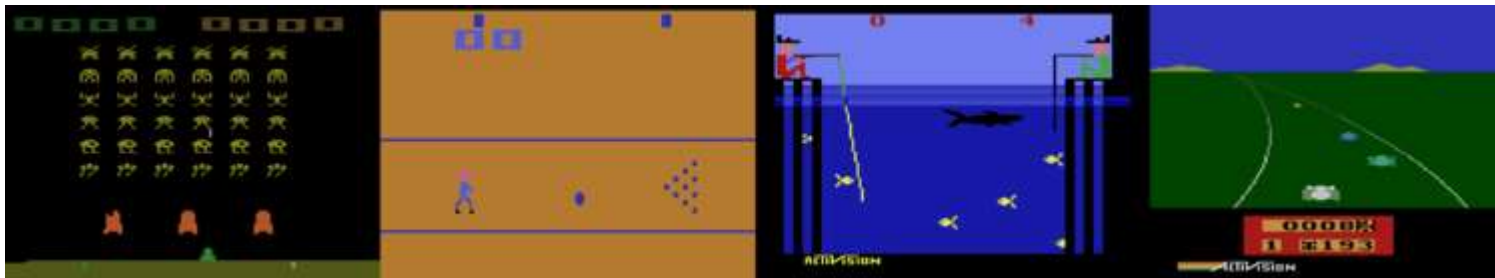
- Computer Go: 3^{361} states



- Inverted pendulum: (x, x', θ, θ')
 - 4-dimensional continuous state space



- Atari: 210x160x3 dimensions (pixel values)



Functions to be Approximated

- Policy: $\delta(s) \rightarrow a$
- Q-function: $Q(s, a) \in \mathfrak{R}$
- Value function: $V(s) \in \mathfrak{R}$

Q-function Approximation

- Let $s = (x_1, x_2, \dots, x_n)^T$

- Linear

$$Q(s, a) \approx \sum_i w_{ai} x_i$$

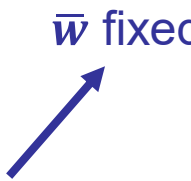
- Non-linear (e.g., neural network)

$$Q(s, a) \approx g(x; w)$$

Gradient Q-learning

- Minimize squared error between Q-value estimate and target
 - Q-value estimate: $Q_{\mathbf{w}}(s, a)$
 - Target: $R(s) + \gamma \max_{a'} Q_{\bar{\mathbf{w}}}(s', a')$

- Squared error:

$$Err(\mathbf{w}) = \frac{1}{2} [Q_{\mathbf{w}}(s, a) - R(s) - \gamma \max_{a'} Q_{\bar{\mathbf{w}}}(s', a')]^2$$


- Gradient

$$\frac{\partial Err}{\partial \mathbf{w}} = [Q_{\mathbf{w}}(s, a) - R(s) - \gamma \max_{a'} Q_{\mathbf{w}}(s', a')] \frac{\partial Q_{\mathbf{w}}(s, a)}{\partial \mathbf{w}}$$

Gradient Q-learning

Initialize weights w at random in $[-1,1]$

Observe current state s

Loop

Select action a and execute it

Receive immediate reward r

Observe new state s'

$$\text{Gradient: } \frac{\partial \text{Err}}{\partial w} = [Q_w(s, a) - r - \gamma \max_{a'} Q_w(s', a')] \frac{\partial Q_w(s, a)}{\partial w}$$

$$\text{Update weights: } w \leftarrow w - \alpha \frac{\partial \text{Err}}{\partial w}$$

$$\text{Update state: } s \leftarrow s'$$

Recap: Convergence of Tabular Q-learning

- Tabular Q-Learning converges to optimal Q-function under the following conditions:

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

- Let $\alpha(s, a) = 1/N(s, a)$
 - Where $N(s, a)$ is # of times that (s, a) is visited

- Q-learning

$$Q(s, a) \leftarrow Q(s, a) + \alpha(s, a) [R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

Convergence of Linear Gradient Q-Learning

- Linear Q-Learning converges under the same conditions:

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

- Let $\alpha_t = 1/t$
- Let $Q_w(s, a) = \sum_i w_i x_i$
- Q-learning

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha_t \left[Q_w(s, a) - r - \gamma \max_{a'} Q_w(s', a') \right] \frac{\partial Q_w(s, a)}{\partial \mathbf{w}}$$

Divergence of non-linear Q-learning

- Even when the following conditions hold

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

non-linear Q-learning may diverge

- Intuition:
 - Adjusting w to increase Q at (s, a) might introduce errors at nearby state-action pairs.

Mitigating divergence

- Two tricks are often used in practice:
 1. Experience replay
 2. Use two networks:
 - Q-network
 - Target network

Experience Replay

- Idea: store previous experiences (s, a, s', r) into a buffer and sample a mini-batch of previous experiences at each step to learn by Q-learning
- Advantages
 - Break correlations between successive updates (more stable learning)
 - Fewer interactions with environment needed to converge (greater data efficiency)

Target Network

- Idea: Use a separate target network that is updated only periodically

repeat for each (s, a, s', r) in mini-batch:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha_t \left[\underbrace{Q_{\mathbf{w}}(s, a)}_{\text{update}} - r - \gamma \max_{a'} \underbrace{Q_{\bar{\mathbf{w}}}(s', a')}_{\text{target}} \right] \frac{\partial Q_{\mathbf{w}}(s, a)}{\partial \mathbf{w}}$$
$$\bar{\mathbf{w}} \leftarrow \mathbf{w}$$

- Advantage: mitigate divergence

Target Network

- Similar to value iteration:

repeat for all s

$$\underbrace{V(s)}_{\text{update}} \leftarrow \max_a R(s) + \gamma \sum_{s'} \Pr(s'|s, a) \underbrace{\bar{V}(s')}_{\text{target}} \quad \forall s$$

$$\bar{V} \leftarrow V$$

repeat for each (s, a, s', r) in mini-batch:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha_t \left[\underbrace{Q_{\mathbf{w}}(s, a)}_{\text{update}} - r - \gamma \max_{a'} \underbrace{Q_{\bar{\mathbf{w}}}(s', a')}_{\text{target}} \right] \frac{\partial Q_{\mathbf{w}}(s, a)}{\partial \mathbf{w}}$$

$$\bar{\mathbf{w}} \leftarrow \mathbf{w}$$

update

target

Deep Q-network

- Google Deep Mind:
- Deep Q-network: Gradient Q-learning with
 - Deep neural networks
 - Experience replay
 - Target network
- Breakthrough: human-level play in many Atari video games

Deep Q-network

Initialize weights w and \bar{w} at random in $[-1,1]$

Observe current state s

Loop

 Select action a and execute it

 Receive immediate reward r

 Observe new state s'

 Add (s, a, s', r) to experience buffer

 Sample mini-batch of experiences from buffer

 For each experience $(\hat{s}, \hat{a}, \hat{s}', \hat{r})$ in mini-batch

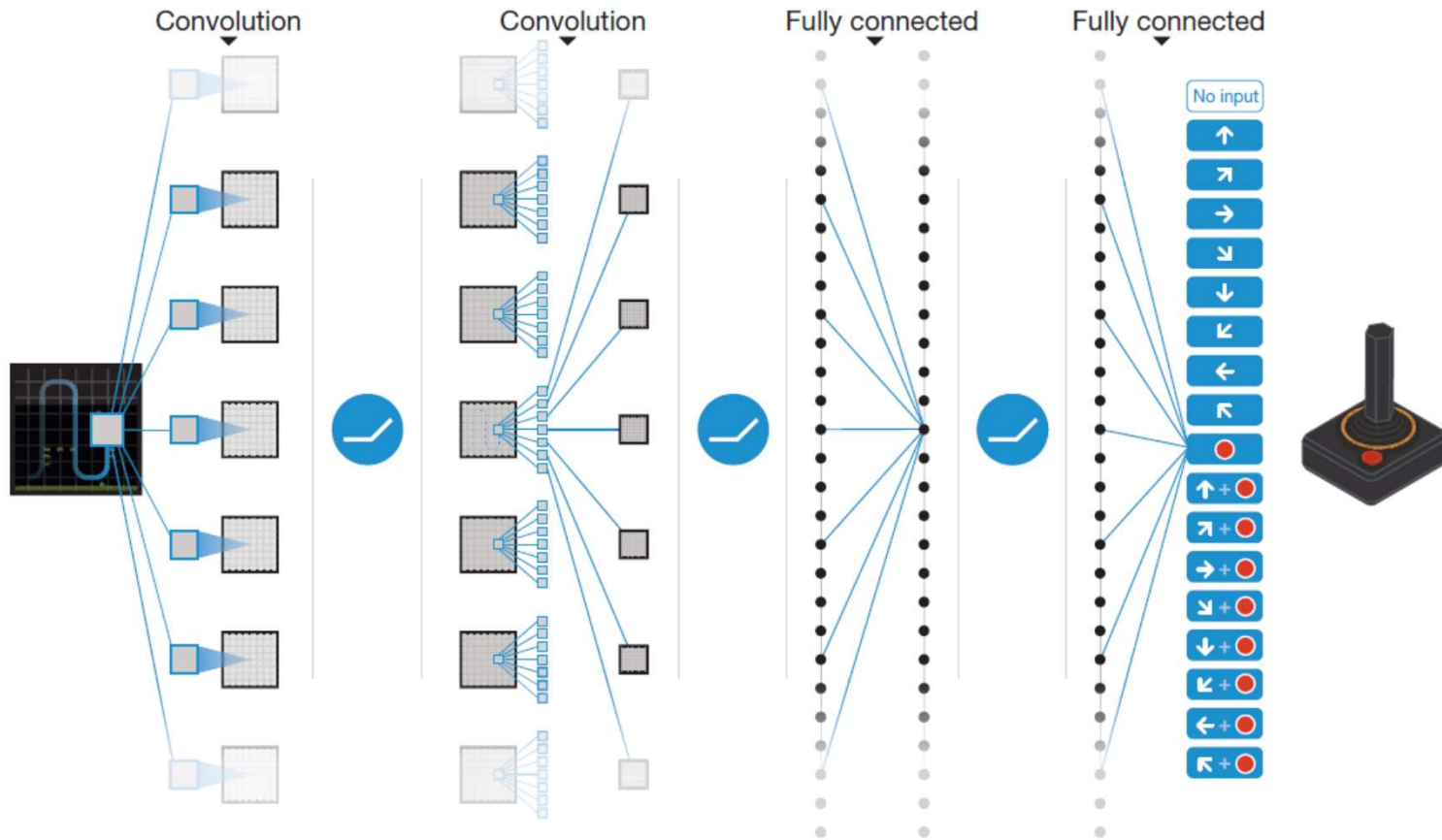
$$\text{Gradient: } \frac{\partial \text{Err}}{\partial w} = [Q_w(\hat{s}, \hat{a}) - \hat{r} - \gamma \max_{\hat{a}'} Q_{\bar{w}}(\hat{s}', \hat{a}')] \frac{\partial Q_w(\hat{s}, \hat{a})}{\partial w}$$

$$\text{Update weights: } w \leftarrow w - \alpha \frac{\partial \text{Err}}{\partial w}$$

Update state: $s \leftarrow s'$

Every c steps, update target: $\bar{w} \leftarrow w$

Deep Q-Network for Atari



DQN versus Linear approx.

