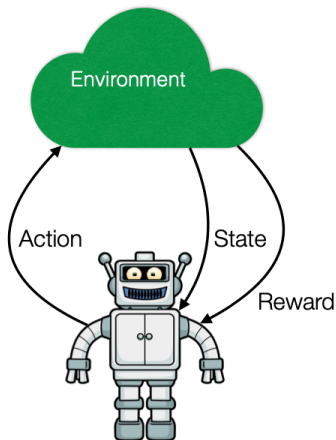# Intro to Deep Q Networks

Jake Grigsby

University of Virginia

April 26, 2020

# Agents and the Environment

Many problems in AI can be thought of as some autonomous agent interacting with an environment:



This concept is formalized as a **Markov Decision Process**...

# Markov Decision Process

> **Definition**
>
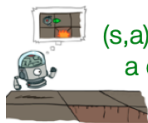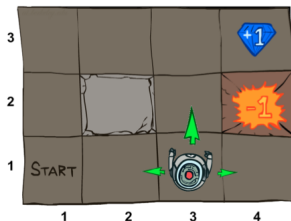> A Markov Decision Process **(MDP)** consists of:
>
> - $\mathcal{S}$, a set of states
> - $\mathcal{A}$, a set of actions
> - $\mathcal{R} \subseteq \mathbb{R}$, a set of rewards
> - a dynamics function $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \to [0, 1]$
>
> $$p(s', r|s, a) = Pr\{\mathcal{S}_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$
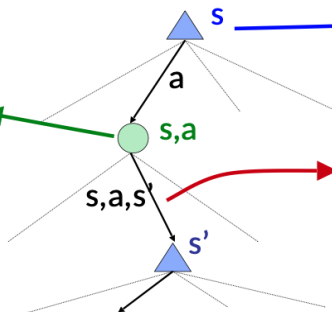
It's common to break the dynamics function $p$ up into a **Transition Function** $T(s, a, s') = \sum_{r \in \mathcal{R}} p(s', r|s, a)$, and a **Reward Function** $R(s, a) = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r|s, a)$

# A Quick Example

# The RL Problem

The goal of RL agents is to find a **policy**[1] $\pi^* : \mathcal{S} \to \mathcal{A}$ that maximizes the *expected discounted return*

$$\pi^* = \underset{\pi}{argmax} \ \underset{\tau \sim \pi}{\mathbb{E}} \left[ \sum_{t=0}^{t=\infty} \gamma^t R_t \right]$$

where $\gamma \in [0, 1)$ is the *discount factor* that lets us deal with non-episodic tasks and $\tau$ is a *trajectory* (a sequence of states and actions that describe the agent's experience)

---

[1]Policies can also be stochastic, in which case they're written $\pi(a|s) : \mathcal{S} \times \mathcal{A} \to [0, 1]$

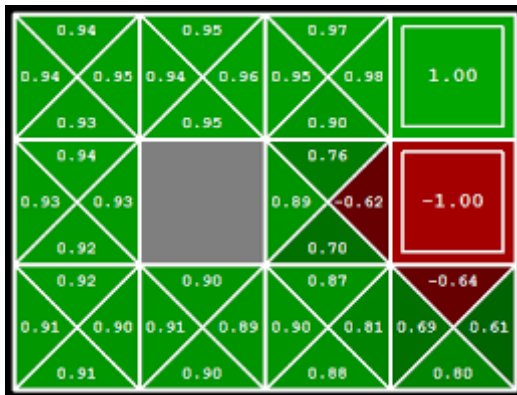# Simplifying Assumptions

We begin by making some assumptions about the task we are trying to solve:

1. The dynamics of the model ($p(s', r|s, a)$) are known
2. $|\mathcal{S}| \ll \infty$
3. $|\mathcal{A}| \ll \infty$

# Generalized Policy Iteration

Solution: Policy Iteration Dynamic Programming



We'll skip these details because knowledge of dynamics is such a limiting assumption in our case. More info can be found in [4]

# Simplifying Assumptions

1. ~~The dynamics of the model ($p(s', r | s, a)$) are known~~
2. $|\mathcal{S}| \ll \infty$
3. $|\mathcal{A}| \ll \infty$

What if the environment dynamics are unknown?

# Value Methods

## Definition

Value methods attempt to learn the optimal *Q Function*

$$Q^*(s, a) = \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{k=0}^{k=\infty} \gamma^k R_{t+k+1} \mid \mathcal{S}_t = s, \mathcal{A}_t = a \right]$$

Why? Because given $Q^*(s, a)$, the optimal policy can easily be computed by

$$\pi^*(s) = \underset{a}{\arg\max} \, Q^*(s, a)$$

# Value Methods: Temporal Difference

- Randomly initialize Q(s, a) and use interactions with the environment as a sample to update this 'bootstrap'
- Updates based on the Bellman Equation:

$$Q^\pi(s, a) = \mathop{\mathbb{E}}_{s'} \left[ r(s, a) + \gamma \mathop{\mathbb{E}}_{a' \sim \pi} \left[ Q^\pi(s', a') \right] \right]$$

# Value Methods: Temporal Difference

- Randomly initialize Q(s, a) and use interactions with the environment as a sample to update this 'bootstrap'

- Updates based on the Bellman Equation:

$$Q^\pi(s, a) = \mathop{\mathbb{E}}_{s'} \left[ r(s, a) + \gamma \mathop{\mathbb{E}}_{a' \sim \pi} \left[ Q^\pi(s', a') \right] \right]$$

- We minimize the *Bellman Error*:

$$(r(s, a) + \gamma \max_{a'} Q(s', a')) - Q(s, a)$$

# Value Methods: Temporal Difference

**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        Take action $A$, observe $R$, $S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma \max_a Q(S', a) - Q(S, A) \big]$
        $S \leftarrow S'$
    until $S$ is terminal

Figure: Q-Learning Pseudo-code [4]

# A Quick Note on Exploration vs. Exploitation

- At each time step, the agent must choose between "exploiting" the action it currently thinks has the best return and "exploring" alternatives to learn more about them.
- Most convergence guarantees assume state *coverage*
  - Every state will be visited an infinite number of times in an infinite number of timesteps.
  - This can be acheived by enforcing:

$$\pi(a|s) > 0, \forall s \in \mathcal{S}$$

# A Quick Note on Exploration vs. Exploitation

- The simplest way to do this is to make an existing policy $\epsilon$-greedy:

$$\pi'(s) = \begin{cases} \pi(s) & \text{with probability } (1 - \epsilon); \\ \pi_{random}(s) & \text{with probability } \epsilon; \end{cases}$$

# A Quick Note on Exploration vs. Exploitation

- The simplest way to do this is to make an existing policy $\epsilon$-greedy:

$$\pi'(s) = \begin{cases} \pi(s) & \text{with probability } (1 - \epsilon); \\ \pi_{random}(s) & \text{with probability } \epsilon; \end{cases}$$

- This can be thought of as injecting noise into the action space
- All of the agent's we'll be talking about use this general approach, but there is a lot of interesting work on motivating agents to explore efficiently. [3] [5]

# Simplifying Assumptions

1. ~~The dynamics of the model ($p(s', r|s, a)$) are known~~
2. ~~$|\mathcal{S}| \ll \infty$~~
3. $|\mathcal{A}| \ll \infty$

What if the state space is too large for dynamic programming?

# Tasks with Large State Spaces

Example: Video Games

- Pixel input makes $|S| = \mathbb{Z}_{256}^{H \times W \times C}$
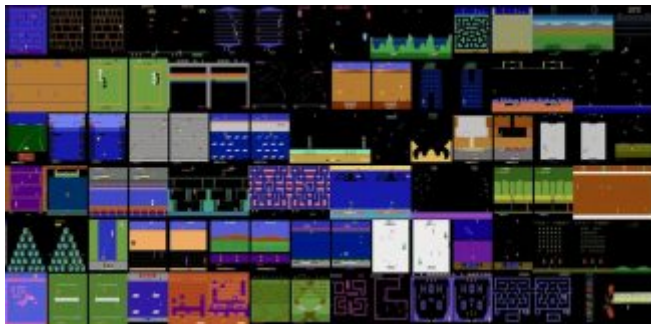- Atari 2600 games make up one of the most popular benchmarks in modern RL.
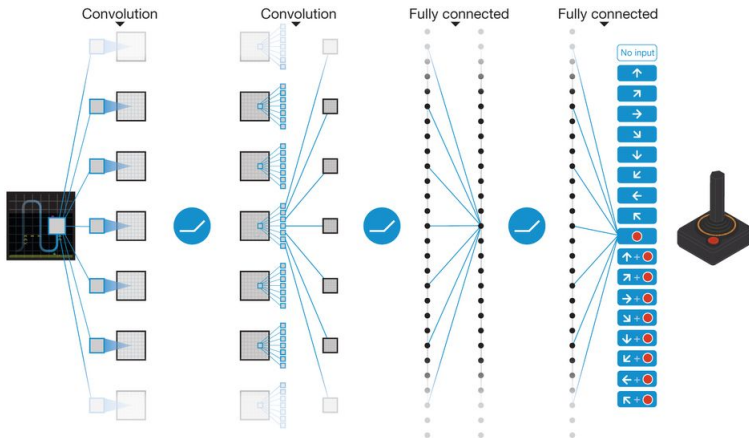


Figure: Games in the Arcade Learning Environment [1] benchmark

# Deep Q Networks (DQN)

- Paramaterize Q with a neural network that can learn to recognize patterns between similar states.

# Deep Q Networks (DQN)

- Train this network to minimize the *Mean Squared Bellman Error* (MSBE)

$$BE(s, a, r, s', d) = (r + \gamma(1 - d) \max_{a'} Q_{\theta'}(s', a')) - Q_\theta(s, a)$$

# Deep Q Networks (DQN)

- Train this network to minimize the *Mean Squared Bellman Error* (MSBE)

$$BE(s, a, r, s', d) = (r + \gamma(1 - d) \max_{a'} Q_{\theta'}(s', a')) - Q_{\theta}(s, a)$$

- Kind of like supervised deep learning!
  - One important difference:
    - The data distribution depends on the parameters (far from i.i.d)

# Deep Q Networks (DQN)

DQNs [2] use a couple tricks to make this more like supervised learning:

1. Create a *replay buffer* $\mathcal{R}$ to store transitions $(s, a, r, s', d)$
   - Randomly sample from this buffer at each training step.
2. Create a *target network* to generate the the bellman error targets.
   - This is a duplicate of the original network that is not trained but is updated with fresh params every $\sim 10000$ steps.

The original DQN was able to learn superhuman policies on many games with dense reward signals!

# Deep Q Networks (DQN)

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory $D$ to capacity $N$

Initialize action-value function $Q$ with random weights $\theta$

Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$

**For** episode = 1, $M$ **do**

    Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

    **For** $t = 1$,T **do**

        With probability $\varepsilon$ select a random action $a_t$

        otherwise select $a_t = \text{argmax}_a Q(\phi(s_t),a;\theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$

        Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$

        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step j+1} \\ r_j + \gamma \max_{a'} \hat{Q}\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$

        Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the network parameters $\theta$
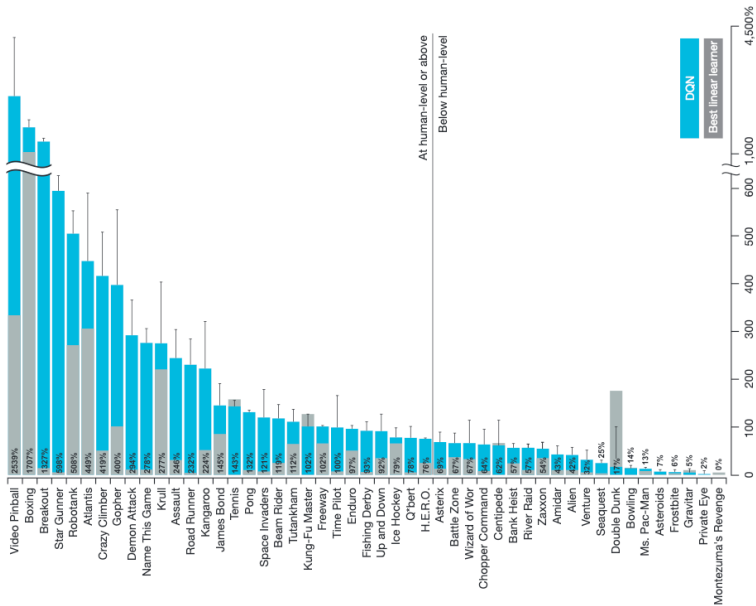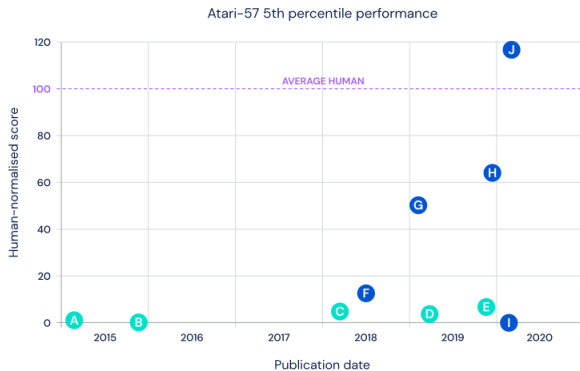
        Every $C$ steps reset $\hat{Q} = Q$

    **End For**

**End For**

Figure: [2]

# Results

# Since 2015...



Atari-57 5th percentile performance

# References I

Marc G Bellemare et al. "The arcade learning environment: An evaluation platform for general agents". In: *Journal of Artificial Intelligence Research* 47 (2013), pp. 253–279.

Volodymyr Mnih et al. "Human-level control through deep reinforcement learning". In: *Nature* 518.7540 (2015), pp. 529–533.

Yuri Burda et al. "Large-scale study of curiosity-driven learning". In: *arXiv preprint arXiv:1808.04355* (2018).

Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

Adrià Puigdomènech Badia et al. "Never Give Up: Learning Directed Exploration Strategies". In: *arXiv preprint arXiv:2002.06038* (2020).