

# Deep Reinforcement Learning

## Seminar: Introduction to Deep Learning

Erik A. Daxberger

February 3, 2017

# Motivation

Ultimate goal of ML/AI research: **general purpose intelligence**, i.e. general intelligent agent, excelling at a wide variety of human-level tasks

Real world is insanely complex → use **video games** as a testbed!

Is it possible to design a single AI agent, being able to play a wide variety of different games, **end-to-end**, at a human level?

# Motivation

Recently, **DeepMind Technologies** published a paper succeeding at that task.



Shortly after, DeepMind was bought by Google for **\$500M** and published a Nature cover paper



# Motivation

(Oh and by the way, they recently managed to beat a human pro at Go)

**So how did they do that?! → deep reinforcement learning!**

## Goal of this talk

- ① good understanding of DeepMind's Atari paper
- ② ability to read other state-of-the-art papers

# Contents

- 1 Introduction
- 2 Reinforcement Learning
- 3 Deep Q-Learning
- 4 Demo: Learning to Play Pong
- 5 Conclusions

# Contents

- 1 Introduction
- 2 Reinforcement Learning
- 3 Deep Q-Learning
- 4 Demo: Learning to Play Pong
- 5 Conclusions

# Machine Learning Paradigms

## ① Supervised learning

- in: example input/output pairs (by a "teacher")
- out: general mapping rule
- example: learn to detect dogs in images → classification

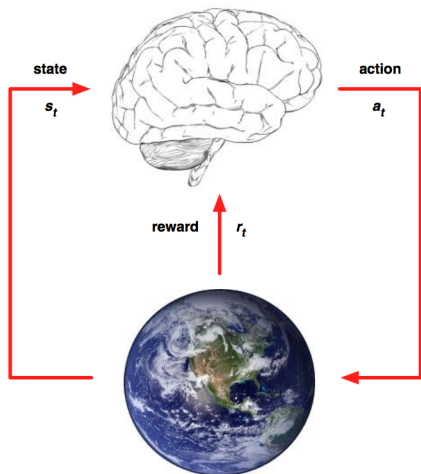
## ② Unsupervised learning

- in: inputs without outputs
- out: hidden structure in the inputs
- example: learn groups of different animals → clustering

## ③ Reinforcement learning

- in: observations and feedback (rewards or punishments) from interacting with an environment
- out: achieve some goal
- example: learn to play Atari games from scratch

# The RL Setting



RL is a **general-purpose framework** for AI

At each timestep  $t$  an **agent** interacts with an **environment**  $\mathcal{E}$  and

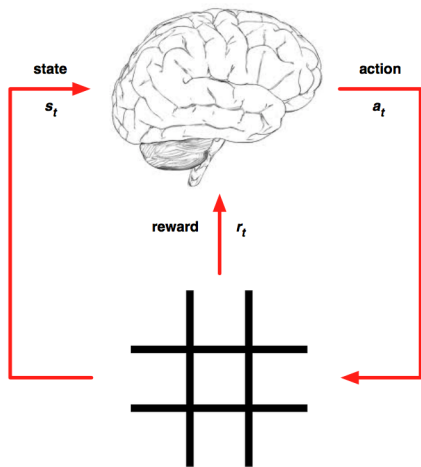
- observes **state**  $s_t \in \mathbb{R}^d$
- receives **reward**  $r_t \in \mathbb{R}$
- executes **action**  $a_t \in A$

Reward signals may be sparse, noisy and delayed.

→ **Agent-environment-loop**



# Example: TicTacToe



$$s_t \in \left\{ \begin{array}{|c|c|c|} \hline o & x & x \\ \hline x & x & o \\ \hline x & o & o \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline x & x & x \\ \hline o & x & o \\ \hline o & o & o \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline o & & \\ \hline x & & o \\ \hline x & & o \\ \hline \end{array}, \dots \right\}$$

$$r_t = \begin{cases} +1 & \text{if we won} \\ -1 & \text{if we lost} \\ 0 & \text{otherwise} \end{cases}$$

$$a_t \in \left\{ \begin{array}{|c|c|c|} \hline x & & \\ \hline & & \\ \hline & & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline & x & \\ \hline & & \\ \hline & & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline & & x \\ \hline & & \\ \hline & & \\ \hline \end{array}, \dots \right\}$$

# Reinforcement Learning vs. Deep Learning

Why don't we just use **Deep Learning**?!

<b>Reinforcement Learning</b>	<b>vs.</b>	<b>Deep Learning</b>
sparse and noisy rewards	vs.	hand-labeled training data
delay between actions and rewards ( $\triangleq$ <i>credit assignment problem</i> )	vs.	direct association between inputs and outputs
highly correlated inputs from a non-stationary data distribution	vs.	i.i.d. data samples

# Fundamental RL Concepts

**Policy**  $\pi$  tells the agent which actions to choose, given states

$$a = \pi(s)$$

**Goal:** Select actions to **maximize future reward**  $\triangleq$  **return**

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$$

**Value function**  $Q^\pi(s, a)$  defines the expected total reward from state  $s$  and action  $a$  under policy  $\pi$  ("How good is action  $a$  in state  $s$ ?" )

$$Q^\pi(s, a) = \mathbb{E}[R_t | s, a] = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s, a]$$

**Optimal value function**

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a)$$

# Example: TicTacToe

$$\pi(s) = \begin{cases} \begin{array}{|c|c|c|} \hline & & x \\ \hline & & \\ \hline & & \\ \hline \end{array} & \text{if } s = \begin{array}{|c|c|c|} \hline x & x & \\ \hline & o & \\ \hline o & & \\ \hline \end{array} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & x & \\ \hline & & \\ \hline \end{array} & \text{if } s = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \\ \dots & \end{cases}$$

Let  $\gamma = 0.9$

$$\begin{aligned} R_t &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \\ &= 0 + 0.9 * 0 + \dots + 0.9^5 * 1 \\ &\approx 0.59 \end{aligned}$$

	<div><div><div></div><div>x</div><div></div></div><div><div></div><div>o</div><div></div></div><div><div></div><div></div><div></div></div></div>	<div><div><div></div><div>o</div><div></div></div><div><div>x</div><div></div><div>x</div></div><div><div></div><div>o</div><div></div></div></div>	...	<div><div><div></div><div>o</div><div></div></div><div><div>o</div><div>x</div><div></div></div><div><div></div><div></div><div>x</div></div></div>
<div><div><div>x</div><div></div><div></div></div><div><div></div><div></div><div></div></div><div><div></div><div></div><div></div></div></div>	0.3	0.5	...	0.9
⋮	⋮	⋮	⋮	⋮
<div><div><div></div><div></div><div></div></div><div><div></div><div></div><div></div></div><div><div></div><div></div><div>x</div></div></div>	0.6	0.1	...	0.8

# Bellman Equation and Value Iteration

Optimal value function can be unrolled recursively  $\rightarrow$  **Bellman equation**

$$\begin{aligned} Q^*(s, a) &= \mathbb{E} [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s, a] \\ &= \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') \middle| s, a \right] \end{aligned}$$

**Value iteration** algorithms estimate  $Q^*$  using iterative **Bellman updates**

$$Q_{i+1}(s, a) = \mathbb{E} \left[ r + \gamma \max_{a'} Q_i(s', a') \middle| s, a \right]$$

This procedure is **guaranteed to converge**, i.e.

$$Q_i \rightarrow Q^* \quad \text{as} \quad i \rightarrow \infty$$

# Function Approximators

These tabular methods are **highly impractical**

- ① all  $Q$ -values are stored separately → **technical challenge**

Example: Chess

- $10^{47}$  states
- 35 possible moves/state
- →  $10^{47} * 35 * 1\text{Byte} \approx 10^{27}$  Zettabytes

(The Internet is 10 Zettabytes)

- ② no generalization over unseen states → **"dumb" approach**

⇒ use a **function approximator** to estimate  $Q$ !

# Contents

- 1 Introduction
- 2 Reinforcement Learning
- 3 Deep Q-Learning**
- 4 Demo: Learning to Play Pong
- 5 Conclusions

## DQN Idea

Approximate  $Q$  by an **ANN** with weights  $\theta \rightarrow$  **deep Q-network (DQN)**

$$Q(s, a; \theta) \approx Q^\pi(s, a)$$

**Loss function** = MSE in Q-Values

$$\mathcal{L}(\theta) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ \left( \underbrace{r + \gamma \max_{a'} Q(s', a'; \theta)}_{\text{target}} - Q(s, a; \theta) \right)^2 \right]$$

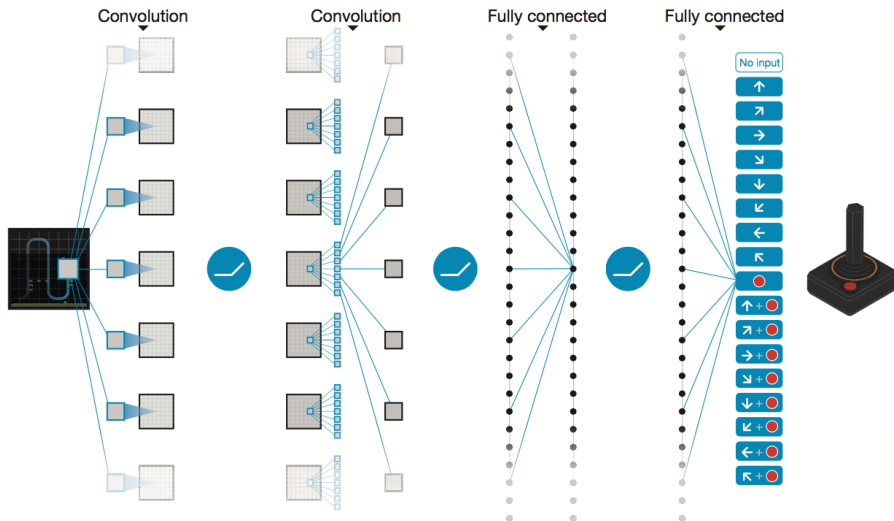
**Q-learning gradient**

$$\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta) - Q(s, a; \theta) \right) \nabla_{\theta} Q(s, a; \theta) \right]$$

$\rightarrow$  Optimize  $\mathcal{L}(\theta)$  with  $\nabla_{\theta} \mathcal{L}(\theta)$  using **stochastic gradient descent**



# DQN Architecture



# Stability Issues & Solutions

Standard Q-learning **oscillates** or **diverges** when using ANNs

- ① Sequential nature of data  $\rightarrow$  non-i.i.d., as successive samples are (highly) correlated
- ② Policy changes rapidly with slight changes to Q-values  $\rightarrow$  policy may oscillate and data distribution may swing between extremes

Tweaks DQN uses to **stabilize training**

- ① **Experience replay**: store agent's experiences  $e_t = (s_t, a_t, r_t, s_{t+1})$  in a data set  $\mathcal{D} = \{e_1, \dots, e_N\}$ , into a *replay memory*  $\rightarrow$  apply minibatch updates to samples  $e \sim \mathcal{D}$
- ② **Frozen target network**: hold previous parameters  $\theta^-$  fixed in the Q-learning target when optimizing  $\mathcal{L}_i(\theta)$ , and update them only periodically, i.e.  $\theta^- \leftarrow \theta$

## DQN Pseudocode

---

**Algorithm:** Deep Q-learning

---

Initialize Q-function with random weights

**for**  $episode = 1, \dots, M$  **do**

**for**  $t = 1, \dots, T$  **do**

        Choose  $a_t = \begin{cases} \text{sampled randomly} & \text{with probability } \epsilon \\ \max_a Q^*(s_t, a; \theta) & \text{otherwise} \end{cases}$

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

        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$

        Sample random minibatch of transitions  $(s_j, a_j, r_j, s_{j+1})$  from  $\mathcal{D}$

        Set  $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{j+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$

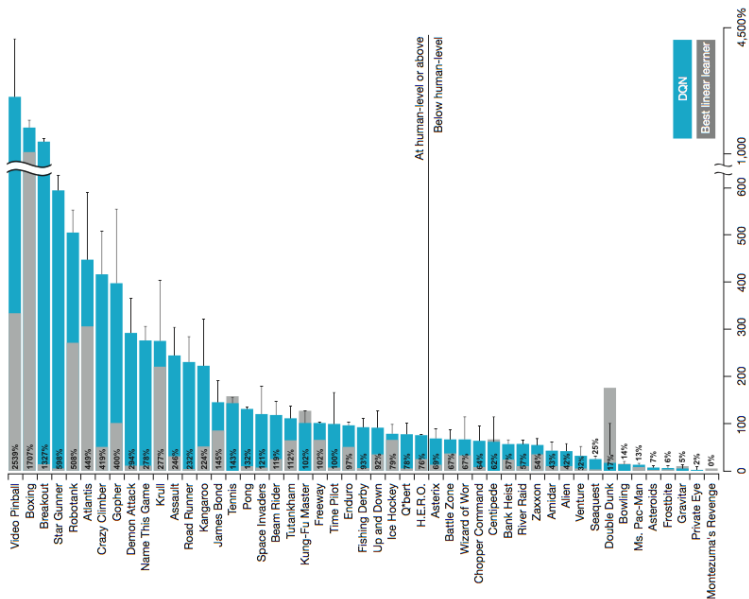
        Perform a gradient descent step on  $(y_j - Q(s_j, a_j; \theta))^2$

**end**

**end**

---

## DQN Atari Results



# Contents

- 1 Introduction
- 2 Reinforcement Learning
- 3 Deep Q-Learning
- 4 Demo: Learning to Play Pong**
- 5 Conclusions

# Framework

**Goal:** Agent that learns how to play Pong, end-to-end

- **input:** raw image frame in form of a  $210 \times 160 \times 3$  byte array, i.e.

$$s_t \in \{0, \dots, 255\}^{100.800}$$

- **output:** distribution over actions  $\rho(a_t)$  ( $\rightarrow$  *stochastic* policy, i.e., we sample  $a_t \sim \rho(a_t)$  in every  $t$ ), with

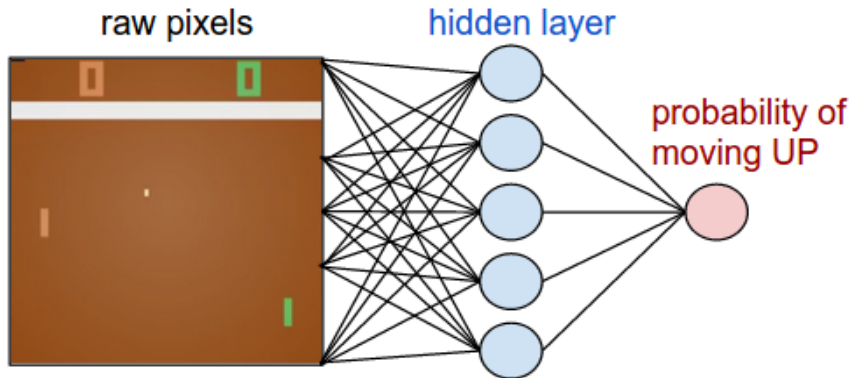
$$a_t \in \{\text{UP}, \text{DOWN}\}$$

- $\rightarrow$  game emulator executes  $a_t$  and emits  $s_{t+1}$  and a **reward**, i.e.,

$$r_t = \begin{cases} +1 & \text{if we score a point} \\ -1 & \text{if our opponent scores a point} \\ 0 & \text{otherwise} \end{cases}$$

# Architecture

**Policy network:** 2-layer fully-connected ANN with 200 hidden nodes



→ **Goal:** find optimal weights  $\theta$  of the policy network

# Implementation

Implemented in Python, using only numpy and **OpenAI Gym**



OpenAI Gym BETA

A toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Go.

```
import gym

env = gym.make("Pong-v0")
s_0 = env.reset()

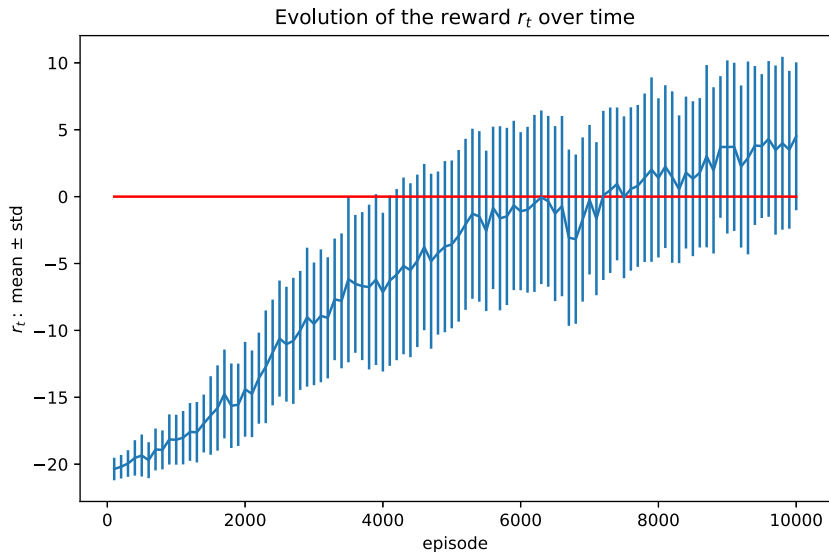
for _ in range(1000):
    env.render()

    a_t = env.action_space.sample()

    s_t1, r_t1, done, info = env.step(a_t)
```



# Results



# Results



# Contents

- 1 Introduction
- 2 Reinforcement Learning
- 3 Deep Q-Learning
- 4 Demo: Learning to Play Pong
- 5 Conclusions**

# Summary

- RL is a **general-purpose AI framework**, modelled through the **agent-environment-loop** of actions, states and rewards
- it can solve problems that are difficult for supervised learning, i.e. involving **sparse, noisy and delayed rewards** and **correlated inputs**
- the goal of an RL agent is to find the **policy maximizing its return**
- this can be done by finding the **optimal value function**, using iterative **Bellman updates**, which is guaranteed to converge
- as storing Q in a table can be demanding for large problems and disables generalization, **function approximators** are commonly used
- for example, an **ANN** can be used to approximate Q, allowing us to optimize Q using standard **deep learning** techniques
- as this process may **oscillate/converge**, DQN uses **experience replay** and a **frozen target network** to stabilize training

## Open Problems and Current Research

Despite the empirical successes, there are still many fundamental **unsolved problems**.

For example, current algorithms

- are **bad at long-term planning** (i.e. act myopically)
- need a **long training time** (i.e., have a low data efficiency)
- are **unable to understand abstract concepts** (i.e. lack a model)

Some **state-of-the-art** papers (Deep RL Workshop @ NIPS 2016)

- Duan, Yan, et al. *RL<sup>2</sup>: Fast Reinforcement Learning via Slow Reinforcement Learning*
- Jaderberg, Max, et al. *Reinforcement Learning with Unsupervised Auxiliary Tasks*
- He, Frank S., et al. *Learning to Play in a Day: Faster Deep Reinforcement Learning by Optimality Tightening*

# Now what about AlphaGo?!

Chess	vs.	Go
$10^{47}$ states	vs.	$10^{170}$ states
rule-based	vs.	intuition-based

So how does AlphaGo work?

**Deep reinforcement learning**  
+  
**Monte Carlo tree search**

→ Silver, David, et al. *Mastering the game of Go with deep neural networks and tree search*. Nature 529.7587 (2016): 484-489.



**Thank you for your attention!**

# References

- Mnih, Volodymyr, et al. *Playing atari with deep reinforcement learning*. arXiv preprint arXiv:1312.5602 (2013).
- Mnih, Volodymyr, et al. *Human-level control through deep reinforcement learning*. Nature 518.7540 (2015): 529-533.
- Sutton, Richard S., and Andrew G. Barto. *Reinforcement learning: An introduction*. Vol. 1. No. 1. Cambridge: MIT press, 1998.
- Karpathy, Andrej. *Deep Reinforcement Learning: Pong from Pixels*. <http://karpathy.github.io/2016/05/31/r1/> (21.01.2017).
- Brockman, Greg, et al. *OpenAI gym*. arXiv preprint arXiv:1606.01540 (2016) and <https://gym.openai.com/> (21.01.2017).
- Silver, David. *Deep Reinforcement Learning*. Tutorial at ICLR (2015), [http://www0.cs.ucl.ac.uk/staff/d.silver/web/Resources\\_files/deep\\_rl.pdf](http://www0.cs.ucl.ac.uk/staff/d.silver/web/Resources_files/deep_rl.pdf) (21.01.2017).