> "The wide world is all about you: you can fence yourselves in, but you cannot forever fence it out."
>
> J. R. R. Tolkien, *The Fellowship of the Ring*

**CMPUT 628**
**Deep RL**

Marlos C. Machado

Class 5 & 6/ 25

Image from THE ONE RING™ Roleplaying Game, Second Edition.

2

# Plan

The general structure of value-based model-free

deep reinforcement learning methods; and DQN.

Marlos C. Machado

# **Please, interrupt me at any time!**

# Reminders & Notes

- You can't leave anymore 😃

- Assignment 1 is due this Saturday, Jan. 24, 2026

- We will release assignment 2 next week, after the ICML deadline (Jan. 28, 2025) ¯\_(ツ)_/¯

- I will release instructions about seminar and paper review during the reading week (Feb 16 – Feb 20)

# Deep RL is about Function Approximation
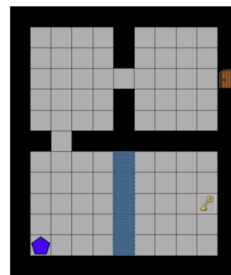
- Why use function approximation?

*"Underparameterization"*
# params < num states



*"Overparameterization"*
# params > num states



Marlos C. Machado

# Deep RL is about Function Approximation

- Why use function approximation?
  - Scalability and generalization

- What is bigger? The number of parameters in the NN or the number of states?

$$\boldsymbol{\theta} \in \mathbb{R}^d \quad q_\pi(s, a; \boldsymbol{\theta}) \approx q_\pi(s, a) \quad d \ll |\mathcal{S}|$$

*"Underparameterization"*
# *params < num states*

**You should try to identify which case you are in**

*"Overparameterization"*
# *params > num states*

– Optimality is impossible, the blanket is too small

– It is about learning a mapping from high-dimensional obs. to the underlying (small) state space.

– We should talk about max. sum of (discounted) rewards

– Often formalized as a Block MDP

Marlos C. Machado

# Observation, Agent State, Environment State

- Atari screens are drawn with 128 colours in a 160 x 210 resolution.
  - There are $128^{210 \times 160}$ ($10^{76220}$) possible images!
  - The Atari 2600 console only has 128 bytes of RAM.

- The number of states in Breakout is estimated  to be between $10^9$ and $10^{11}$
  - If we use 32 bits to represent a state, it could take up to 400 GB!
  - But for comparison, the number of states Chess and Go are estimated to have is $10^{46}$  and $10^{172}$, respectively.

- The observation is not a state. There are two types of state: the environment state and the agent state.
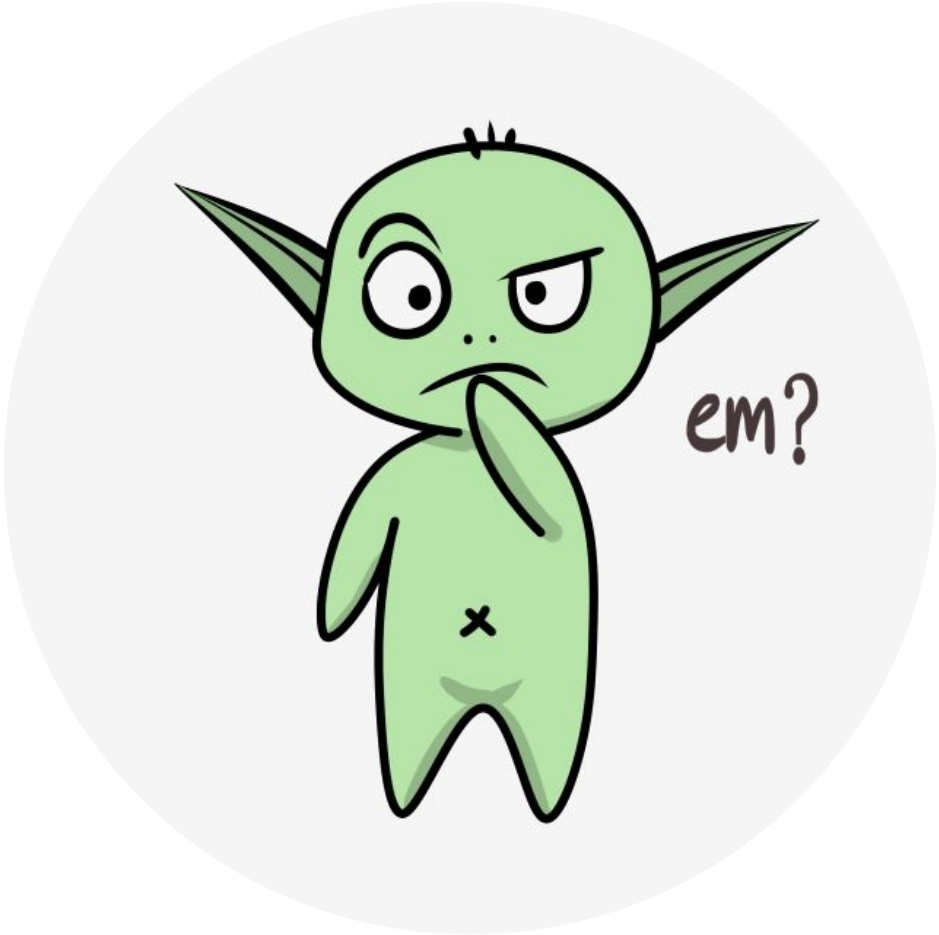
# Function Approximation

- In these next classes we'll be approximating value functions, later we will talk about approximating models or directly approximating policies.

- We are studying deep reinforcement learning because of its scalability and its ability (or promise) to learn representations.

- It is easy to take for granted the representation learning aspect, this is why your assignment 2 is partly about designing features :-)
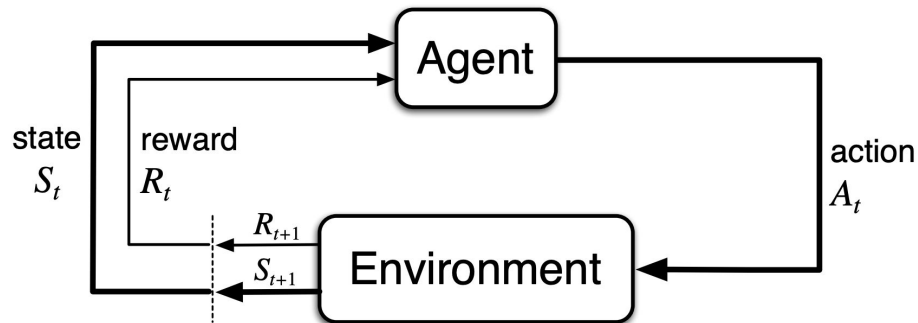
# Example: Jumping Task [Tachet des Combes et al., 2018]



**Which features should you use?**

Marlos C. Machado

em?

# The Main Components of Value-based Model-free Methods

- At least for now, deep RL algorithms are more than RL algorithms with NNs

- In deep RL, the agent-environment interaction *is the same* as in traditional RL, but deep RL agents tend to be quite different from RL agents.

  - They tend to have many more components, they are more complicated

## … but before that…

# In traditional RL, how would we use NNs for funct. approx.?



$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \Big[ R_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q(O_{t+1}, a'; \boldsymbol{\theta}_t) - Q(O_t, A_t; \boldsymbol{\theta}_t) \Big] \nabla_{\boldsymbol{\theta}_t} Q(O_t, A_t; \boldsymbol{\theta}_t)$$

**Why don't we use this?**

Marlos C. Machado

# Online Q-Learning with Neural Networks (OQLNN)

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \left[ R_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q(O_{t+1}, a'; \boldsymbol{\theta}_t) - Q(O_t, A_t; \boldsymbol{\theta}_t) \right] \nabla_{\boldsymbol{\theta}_t} Q(O_t, A_t; \boldsymbol{\theta}_t)$$

- It can be quite unstable

  - Consecutive samples end up being quite correlated

  - Generalization between the current and the next state can be quite dangerous for bootstrapping

- We don't *have to* process one sample at a time

  - We can't parallelize learning with one sample at a time

  - We might want to use a sample more than once

Marlos C. Machado

# The Deadly Triad

- Instabilities arise when you combine:
  - Function approximation,
  - Bootstrapping, and
  - Off-policy learning.

- This happens even for prediction with linear function approximation.

- We are far from having theoretical results really justifying deep RL, so the vast majority of claims I'll make here will be based on empirical data.

For a manuscript with an actual empirical analysis about this issue in deep RL, see reference below.

*Hado van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, Joseph Modayil: Deep Reinforcement Learning and the Deadly Triad. CoRR abs/1812.02648 (2018)*

em?

# Deep Reinforcement Learning

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



**Data sampling strategy**

**Data transformation**

**Backbone of the neural network**

Obj.

Objective func. reward max.

Fixed size buffer to store transitions (often FIFO) [Lin et al., 1991]

Observation augmentation [Tao et al., 2023]

Aux. obj. func. (aux. tasks) [Jaderberg et al., 2017]

Agent-state

# Structural Outline of Model-Free Value-based Deep RL Algs.



**Data sampling strategy**
- Can decorrelate samples
- Allows for parallelism (mini-batch training) and for samples being used more than once
- Often requires off-policy learning
- Often, most recent sample is not processed right away

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



Data transformation

- Despite common beliefs, we rarely feed "raw" inputs to the NN

- For images, people sometimes grayscale, downsample, crop

- For proprioception, people can normalize, apply the Fourier transform, and so on

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



Auxiliary Inputs
[Tao et al., 2023]

- In most of the high-profile results, the NN received more than the *last* obs. transformed

- Examples include: last *k* frames, uncertainty estimates, predictions

Marlos C. Machado

22

# Structural Outline of Model-Free Value-based Deep RL Algs.



**Neural Network**

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



**NN Backbone**

- First couple of layers that will lead to the learned features, $\phi$

- Examples include MLPs, CNNs, and ResNets

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.

Obj.

$f$

Input

Neural
network

$\varphi$

Aux.
Inputs

Experience Buffer

Representation

Aux.
Obj.

- Still the NN, but here we depict the representation

- It can be seen as the agent-state

- It can be recurrent, but training gets trickier because of the order in which samples are processed

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



**Main objective function:**

- The loss/obj that induces reward maximization (value estimation), used for action selection

- It is often some TD-like obj. func.

- I'm conflating loss and obj. here

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



Obj.

$f$

Input

Neural
network

$\varphi$

Aux.
Inputs

Experience Buffer

Auxiliary tasks
[Jaderberg et al., 2017]

Aux.
Obj.

- Predicting more than the value estimate can be quite helpful when training the NN (GVFs)

- Changes the opt. landscape

- Learn better representations

Marlos C. Machado

# Structural Outline of Model-Free Value-based Deep RL Algs.



**Different algorithms are different instantiations of these boxes.**
**This will be the major bulk of our course!**

Marlos C. Machado

em?

# Why use games as an evaluation platform?

*"And some things that should not have been forgotten were lost. History became legend. Legend became myth …."*

*— Galadriel in The Lord of the Rings: The Fellowship of the Ring, The Lord of the Rings*

# Why do we use games as an evaluation platform?

- Games have many useful properties for scientific experimentation:
    - They are fully controllable and have well-defined rules
    - They are free of experimenter's bias
    - They are relatable and challenging
    - They have well-defined metrics of success
- Games require increasingly complex sets of skills that are generally useful:
    - Colours, numbers, pattern matching
    - Effects of actions, short-term planning, long-term planning
    - Strategic thinking, problem solving
    - Cooperation, social skills
- Games are convenient because experts already developed them for us

Marlos C. Machado

Influenced by many discussions with Michael Bowlin

em?

# Atari 2600 Games

em?

35

# Deep Q-Networks (DQN)

Marlos C. Machado

# The beginning of it all: Mnih et al. (2013)

**Playing Atari with Deep Reinforcement Learning**

Volodymyr Mnih    Koray Kavukcuoglu    David Silver    Alex Graves    Ioannis Antonoglou

Daan Wierstra    Martin Riedmiller

DeepMind Technologies

{vlad,koray,david,alex.graves,ioannis,daan,martin.riedmiller} @ deepmind.com

**Abstract**

We present the first deep learning model to successfully learn control policies di-
rectly from high-dimensional sensory input using reinforcement learning. The
model is a convolutional neural network, trained with a variant of Q-learning,
whose input is raw pixels and whose output is a value function estimating future
rewards. We apply our method to seven Atari 2600 games from the Arcade Learn-
ing Environment, with no adjustment of the architecture or learning algorithm. We
find that it outperforms all previous approaches on six of the games and surpasses
a human expert on three of them.

The same neural network architecture, input space, hyperparameters, and, to some degree, action space, were shared across all games

| | B. Rider | Breakout | Enduro | Pong | Q*bert | Seaquest | S. Invaders |
|---|---|---|---|---|---|---|---|
| **Random** | 354 | 1.2 | 0 | −20.4 | 157 | 110 | 179 |
| **Sarsa** [3] | 996 | 5.2 | 129 | −19 | 614 | 665 | 271 |
| **Contingency** [4] | 1743 | 6 | 159 | −17 | 960 | 723 | 268 |
| **DQN** | **4092** | **168** | **470** | **20** | **1952** | **1705** | **581** |
| **Human** | 7456 | 31 | 368 | −3 | 18900 | 28010 | 3690 |
| **HNeat Best** [8] | 3616 | 52 | 106 | 19 | 1800 | 920 | **1720** |
| **HNeat Pixel** [8] | 1332 | 4 | 91 | −16 | 1325 | 800 | 1145 |
| **DQN Best** | **5184** | **225** | **661** | **21** | **4500** | **1740** | 1075 |

Table 1: The upper table compares average total reward for various learning methods by running
an $\epsilon$-greedy policy with $\epsilon = 0.05$ for a fixed number of steps. The lower table reports results of
the single best performing episode for HNeat and DQN. HNeat produces deterministic policies that
always get the same score while DQN used an $\epsilon$-greedy policy with $\epsilon = 0.05$.

# The beginning of it all: Mnih et al. (2013)

The same neural network

**Startups**

## Google Acquires Artificial Intelligence Startup DeepMind For More Than $500M

Join TechCrunch+

Login

Search

Catherine Shu  @catherineshu / 6:20 PM MST • January 26, 2014

Comment

We present the first deep learning model to successfully learn control policies directly from high-dimensional sensory input using reinforcement learning. The model is a convolutional neural network, trained with a variant of Q-learning, whose input is raw pixels and whose output is a value function estimating future rewards. We apply our method to seven Atari 2600 games from the Arcade Learning Environment, with no adjustment of the architecture or learning algorithm. We find that it outperforms all previous approaches on six of the games and surpasses a human expert on three of them.

| | | | | | | | nvaders |
|---|---|---|---|---|---|---|---|
| | | | | | | | 179 |
| | | | | | | | 271 |
| Contingency [4] | 1743 | 6 | 159 | −17 | 960 | 723 | 268 |
| DQN | 4092 | 168 | 470 | 20 | 1952 | 1705 | 581 |
| Human | 7456 | 31 | 368 | −3 | 18900 | 28010 | 3690 |
| HNeat Best [8] | 3616 | 52 | 106 | 19 | 1800 | 920 | 1720 |
| HNeat Pixel [8] | 1332 | 4 | 91 | −16 | 1325 | 800 | 1145 |
| DQN Best | 5184 | 225 | 661 | 21 | 4500 | 1740 | 1075 |

Table 1: The upper table compares average total reward for various learning methods by running an $\epsilon$-greedy policy with $\epsilon = 0.05$ for a fixed number of steps. The lower table reports results of the single best performing episode for HNeat and DQN. HNeat produces deterministic policies that always get the same score while DQN used an $\epsilon$-greedy policy with $\epsilon = 0.05$.

Marlos C. Machado

# … and finally

# And the results caught everyone's attention



Marlos C. Machado

# DQN: A first complete instantiation of our structural outline

Transitions are sampled randomly, so the agent can update with multiple samples in parallel

i.i.d.?

Stores transitions⟨s, a, s', r⟩ into a FIFO queue. They were generated by many policies.

Resizing
Cropping
Luminance

Input

Experience Buffer

Aux.
Inputs

Frame
stacking

CNN

φ

Obj.

Buffer size: 1M
Minibatch size: 32
Update frequency: 4

em?

# DQN objective function and loss function I

$$\mathcal{L}^{\mathrm{DQN}} = \mathbb{E}_{(o,a,r,o') \sim U(\mathcal{D})} \left[ \left( R_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q(O_{t+1}, a'; \boldsymbol{\theta}^-) - Q(O_t, A_t; \boldsymbol{\theta}_t) \right)^2 \right]$$

Expectation is over the uniform
distribution of sampled transitions

Target network

$$Y(R_{t+1}, O_{t+1}; \boldsymbol{\theta}^-) = R_{t+1} + \gamma \max_{a' \in \mathcal{A}} Q(O_{t+1}, a'; \boldsymbol{\theta}^-)$$

$$\mathcal{L}^{\mathrm{DQN}} = \mathbb{E}_{(o,a,r,o') \sim U(\mathcal{D})} \left[ \left( Y(R_{t+1}, O_{t+1}; \boldsymbol{\theta}^-) - Q(O_t, A_t; \boldsymbol{\theta}_t) \right)^2 \right].$$

Regression

Marlos C. Machado

# DQN objective function and loss function II

- DQN, at least Nature-DQN, had way more than an experience replay buffer and a target network

  - To use the same hyperparameters across all the games, the TD error and the loss function had to be in the same scale:
    - Clip the rewards to be [-1, 0, 1]
    - Clip the TD error be between −1 and 1 in the gradient update
    - RMSProp (Adam is now more common)

  - The experience replay buffer needs to be partially filled before training can occur
    - Uniform random policy for 50k frames
    - ε-greedy policy, with ε decaying from 1.0 to 0.1 within the first 1M frames

Marlos C. Machado

# More design choices for DQN

- Observations: Four stacked frames in an attempt to make it Markovian
  I call the additional three frames *auxiliary inputs* (Tao et al., 2023)
  - Instead of 210 × 160 pixels with 128 colours, the observation is pre-processed. It is rescaled and recoloured to an image of size 84 × 84 where the luminance values are extracted from the original image and used to encode each pixel (grayscaling).

- Neural network architecture
  - 3 convolutional layers, a hidden layer, and an output layer
    Conv. layers use 32 (8 × 8), 64 (4 × 4), and 64 (3 × 3)
    filters with stride 4, 2, and 1 respectively. These layers are
    flattened (3,136 units) to connect to a hidden layer of
    512 units, which is connected to the output layer.
    The output layer size ranges from 4 to 18 units.

  - ReLus everywhere.

# DQN and Neural Fitted-Q Iteration

- Neural Fitted Q-Iteration (NFQI) is DQN's predecessor, *they are not the same*.

- Neural Fitted Q-Iteration (NFQI):

  ○ It trains a neural network from scratch in each iteration.

  ○ It keeps all the data it has seen around (the most natural instantiation of the algorithm collects all data beforehand). Thus the intermediate Q-values do not impact the transitions observed by the agent (although the authors do mention a variant in which the dataset is augmented with more data).

  ○ It has some sort of target network as the target values the neural network is regressing to are fixed to the values of the previous iteration, implying the frequency in which the target network is updated is at least the number of samples in the collected dataset.

  ○ It was evaluated only in classic simple tasks such as Mountain Car and Pole Balancing.

# Target networks and experience replay buffers aren't *hacks*

- Replay buffers give us more i.i.d. samples, sample reuse and parallelization

    ○ It can be seen as a model. Dyna was always praised, why the issue with a replay buffer?

    ○ The universe doesn't shuffle data, though

    ○ I find it weird to not perform updates with the sample the agent just saw (but it often doesn't help)

    ○ The size of the replay buffer is a key (and very sensitive) hyperparameter in deep RL agents

- Target networks approximate the fixed targets in supervised learning

    ○ Neural Fitted-Q Iteration and other methods did that before

    ○ But they can slow down learning because the agent is always regressing towards a stale value

    ○ The rate at which the target network is updated is a rather sensitive hyperparameter

    Replay ratio: the number of environment steps taken per gradient step.

Marlos C. Machado

em?

# DQN's evaluation methodology

- Final performance was reported after training an agent with 200 million frames

  - Each action was repeated by the agent 4 times, this parameter is known as *frame skip*

  - Four actions were selected by the agent between successive updates

- Episode terminated upon loss of a life

- Minimal action set was used

> - 200M frames
> - 50M action selections
> - 12.5M grad. updates

- There was an evaluation phase in which the performance of the *best* checkpoint obtained during learning was evaluated 30 times. Learning happened once

- The ALE used to be deterministic, so they used 0-30 no-ops to provide *some* randomization

Marlos C. Machado

# On the comparison between different papers

| Game | Random Play | Best Linear Learner | Contingency (SARSA) | Human | DQN (± std) | Normalized DQN (% Human) |
|------|-------------|---------------------|---------------------|-------|-------------|--------------------------|
| Alien | 227.8 | 939.2 | 103.2 | 6875 | 3069 (±1093) | 42.7% |
| Amidar | 5.8 | 103.4 | 183.6 | 1676 | 739.5 (±3024) | 43.9% |
| Assault | 222.4 | 628 | 537 | 1496 | 3359(±775) | 246.2% |
| Asterix | 210 | 987.3 | 1332 | 8503 | 6012 (±1744) | 70.0% |
| Asteroids | 719.1 | 907.3 | 89 | 13157 | 1629 (±542) | 7.3% |
| Atlantis | 12850 | 62687 | 852.9 | 29028 | 85641(±17600) | 449.9% |
| Bank Heist | 14.2 | 190.8 | 67.4 | 734.4 | 429.7 (±650) | 57.7% |
| Battle Zone | 2360 | 15820 | 16.2 | 37800 | 26300 (±7725) | 67.6% |
| Beam Rider | 363.9 | 929.4 | 1743 | 5775 | 6846 (±1619) | 119.8% |
| Bowling | 23.1 | 43.9 | 36.4 | 154.8 | 42.4 (±88) | 14.7% |
| Boxing | 0.1 | 44 | 9.8 | 4.3 | 71.8 (±8.4) | 1707.9% |
| Breakout | 1.7 | 5.2 | 6.1 | 31.8 | 401.2 (±26.9) | 1327.2% |
| Centipede | 2091 | 8803 | 4647 | 11963 | 8309(±5237) | 63.0% |
| Chopper Command | 811 | 1582 | 16.9 | 9882 | 6687 (±2916) | 64.8% |
| Crazy Climber | 10781 | 23411 | 149.8 | 35411 | 114103 (±22797) | 419.5% |
| Demon Attack | 152.1 | 520.5 | 0 | 3401 | 9711 (±2406) | 294.2% |
| Double Dunk | -18.6 | -13.1 | -16 | -15.5 | -18.1 (±2.6) | 17.1% |
| Enduro | 0 | 129.1 | 159.4 | 309.6 | 301.8 (±24.6) | 97.5% |
| Fishing Derby | -91.7 | -89.5 | -85.1 | 5.5 | -0.8 (±19.0) | 93.5% |
| Freeway | 0 | 19.1 | 19.7 | 29.6 | 30.3 (±0.7) | 102.4% |
| Frostbite | 65.2 | 216.9 | 180.9 | 4335 | 328.3 (±250.5) | 6.2% |
| Gopher | 257.6 | 1288 | 2368 | 2321 | 8520 (±3279) | 400.4% |
| Gravitar | 173 | 387.7 | 429 | 2672 | 306.7 (±223.9) | 5.3% |
| H.E.R.O. | 1027 | 6459 | 7295 | 25763 | 19950 (±158) | 76.5% |
| Ice Hockey | -11.2 | -9.5 | -3.2 | 0.9 | -1.6 (±2.5) | 79.3% |

We always report the game score, because that's what we care about. We often set $\gamma < 1$, though. It is a better surrogate objective due to instabilities when setting $\gamma = 1$.

# Standardizing experimentation and introducing stochasticity

## Revisiting the Arcade Learning Environment:
## Evaluation Protocols and Open Problems for General Agents

**Marlos C. Machado**                                    MACHADO@UALBERTA.CA
*University of Alberta, Edmonton, Canada*

**Marc G. Bellemare**                                    BELLEMARE@GOOGLE.COM
*Google Brain, Montréal, Canada*

**Erik Talvitie**                                        ERIK.TALVITIE@FANDM.EDU
*Franklin & Marshall College, Lancaster, USA*

**Joel Veness**                                          AIXI@GOOGLE.COM
*DeepMind, London, United Kingdom*

**Matthew Hausknecht**                          MATTHEW.HAUSKNECHT@MICROSOFT.COM
*Microsoft Research, Redmond, USA*

**Michael Bowling**                                      MBOWLING@UALBERTA.CA
*University of Alberta, Edmonton, Canada*
*DeepMind, Edmonton, Canada*

### Abstract

The Arcade Learning Environment (ALE) is an evaluation platform that poses the challenge of building AI agents with general competency across dozens of Atari 2600 games. It supports a variety of different problem settings and it has been receiving increasing attention from the scientific community, leading to some high-profile success stories such as the much publicized Deep Q-Networks (DQN). In this article we take a big picture look at how the ALE is being used by the research community. We show how diverse the evaluation methodologies in the ALE have become with time, and highlight some key concerns when evaluating agents in the ALE. We use this discussion to present some methodological best practices and provide new benchmark results using these best practices. To further the progress in the field, we introduce a new version of the ALE that supports multiple game

# Which one is better?



Marlos C. Machado

# It does happen!



Sarsa ($\lambda$)

DQN

# Determinism in the Arcade Learning Environment

- ## The ALE was deterministic:
  - all episodes have the same start state and,
  - the same sequence of actions will always lead to the same outcome.

55

# The Brute [Bellemare et al., 2015]



Marlos C. Machado

# Performance in the deterministic ALE

# Stochasticity model – *Sticky actions*

$$A_t = \begin{cases} a, & \text{with prob.} \quad 1 - \varsigma, \\ a_{t-1}, & \text{with prob.} \quad \varsigma. \end{cases}$$

**Input**



**Execution**

# Performance in the stochastic ALE

# Different RL methods were comparable again. We did not even emphasize benchmarking too much (it didn't matter)

Machado, Bellemare, Talvitie, Veness, Hausknecht, & Bowling

Table 8: Sarsa(λ) + Blob-PROST results across 60 games. See Appendix B for details.

554

Revisiting the ALE: Evaluation Protocols and Open Problems

Table 9: DQN results across 60 games. See Appendix B for details.

555

Marlos C. Machado

em?

# We can always update our beliefs

# Next class

- ● **What I plan to do:**

  - ○ Value-based Model-free Methods Objective Functions: Double Learning

- ● **What I recommend YOU to do for next class:**

  - ○ Read the lecture notes and the Double DQN paper.

    *van Hasselt, H., Guez, A., and Silver, D. (2016). Deep reinforcement learning with double Q-learning. In Proceedings of the Conference on Artificial Intelligence, pages 2094–2100. Preprint made available on September 22, 2015.*

  - ○ Start thinking about Assignment 2!