

Course material at:

<https://github.com/lyeskhalil/mlbootcamp2022/>

# Advanced Topics

Friday

FASE ML Bootcamp

Based on material from various sources:

Data augmentation, Hyperparameter tuning, CNN understanding, Attention, GANs:

<http://cs231n.stanford.edu/schedule.html>

Uncertainty in Deep Learning: <https://www.gatsby.ucl.ac.uk/~balaji/balaji-uncertainty-talk-cifar-dlrl.pdf>

Graph Neural Networks: <http://tkipf.github.io/misc/SlidesCambridge.pdf>

Reinforcement Learning: [https://www.davidsilver.uk/wp-content/uploads/2020/03/deep\\_rl\\_tutorial\\_small\\_compressed.pdf](https://www.davidsilver.uk/wp-content/uploads/2020/03/deep_rl_tutorial_small_compressed.pdf)

**2022 Free PDF book on all of ML by Kevin Murphy (Google):** <https://probml.github.io/pml-book/book1.html>

# Tutorial: Deep Reinforcement Learning

David Silver, Google DeepMind

# Reinforcement Learning in a nutshell

RL is a general-purpose framework for decision-making

- ▶ RL is for an **agent** with the capacity to **act**
- ▶ Each **action** influences the agent's future **state**
- ▶ Success is measured by a scalar **reward** signal
- ▶ Goal: **select actions to maximise future reward**

# Deep Learning in a nutshell

DL is a general-purpose framework for representation learning

- ▶ Given an **objective**
- ▶ Learn **representation** that is required to achieve objective
- ▶ Directly from **raw inputs**
- ▶ Using minimal domain knowledge

# Deep Reinforcement Learning: $AI = RL + DL$

We seek a single agent which can solve any human-level task

- ▶ RL defines the objective
- ▶ DL gives the mechanism
- ▶  $RL + DL =$  general intelligence

# Examples of Deep RL @DeepMind

- ▶ **Play** games: Atari, poker, Go, ...
- ▶ **Explore** worlds: 3D worlds, Labyrinth, ...
- ▶ **Control** physical systems: manipulate, walk, swim, ...
- ▶ **Interact** with users: recommend, optimise, personalise, ...

# Outline

Introduction to Deep Learning

Introduction to Reinforcement Learning

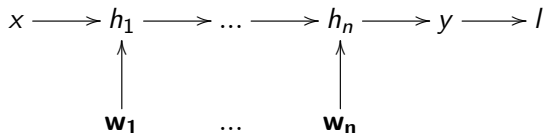
Value-Based Deep RL

Policy-Based Deep RL

Model-Based Deep RL

# Deep Representations

- ▶ A **deep representation** is a composition of many functions



- ▶ Its gradient can be **backpropagated** by the chain rule

$$\begin{array}{ccccccc} \frac{\partial l}{\partial x} & \xleftarrow{\frac{\partial h_1}{\partial x}} & \frac{\partial l}{\partial h_1} & \xleftarrow{\frac{\partial h_2}{\partial h_1}} & \dots & \xleftarrow{\frac{\partial h_n}{\partial h_{n-1}}} & \frac{\partial l}{\partial h_n} & \xleftarrow{\frac{\partial y}{\partial h_n}} & \frac{\partial l}{\partial y} \\ & & \downarrow \frac{\partial h_1}{\partial w_1} & & & & \downarrow \frac{\partial h_n}{\partial w_n} & & \\ & & \frac{\partial l}{\partial \mathbf{w}_1} & & \dots & & \frac{\partial l}{\partial \mathbf{w}_n} & & \end{array}$$



# Deep Neural Network

A **deep neural network** is typically composed of:

- ▶ Linear transformations

$$h_{k+1} = Wh_k$$

- ▶ Non-linear activation functions

$$h_{k+2} = f(h_{k+1})$$

- ▶ A loss function on the output, e.g.
  - ▶ Mean-squared error  $l = ||y^* - y||^2$
  - ▶ Log likelihood  $l = \log \mathbb{P}[y^*]$

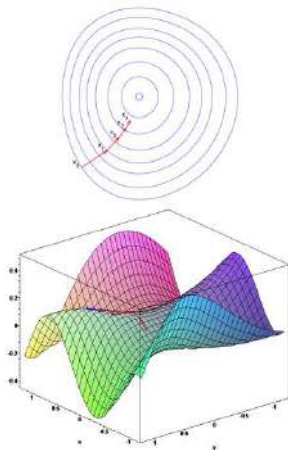
# Training Neural Networks by Stochastic Gradient Descent

- ▶ Sample gradient of expected loss  $L(\mathbf{w}) = \mathbb{E}[l]$

$$\frac{\partial l}{\partial \mathbf{w}} \sim \mathbb{E} \left[ \frac{\partial l}{\partial \mathbf{w}} \right] = \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$

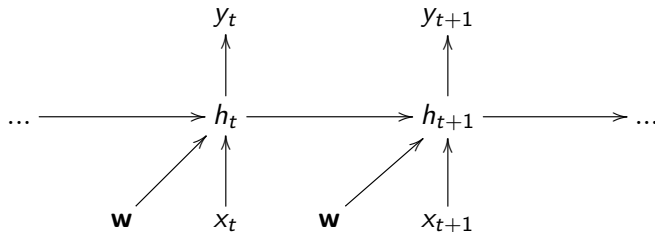
- ▶ Adjust  $\mathbf{w}$  down the sampled gradient

$$\Delta \mathbf{w} \propto \frac{\partial l}{\partial \mathbf{w}}$$

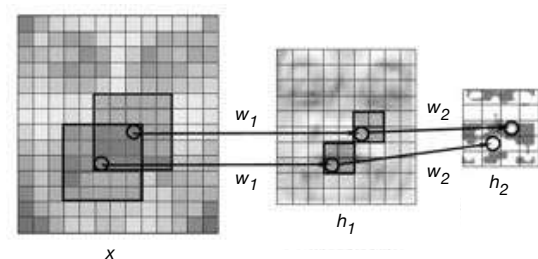


# Weight Sharing

Recurrent neural network shares weights between time-steps



Convolutional neural network shares weights between local regions



# Outline

Introduction to Deep Learning

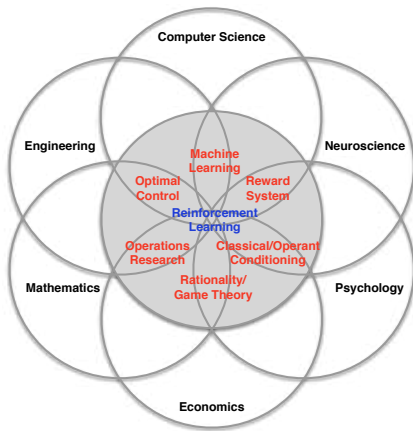
Introduction to Reinforcement Learning

Value-Based Deep RL

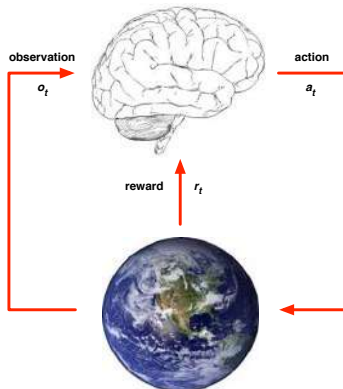
Policy-Based Deep RL

Model-Based Deep RL

# Many Faces of Reinforcement Learning



# Agent and Environment



- ▶ At each step  $t$  the agent:
  - ▶ Executes action  $a_t$
  - ▶ Receives observation  $o_t$
  - ▶ Receives scalar reward  $r_t$
- ▶ The environment:
  - ▶ Receives action  $a_t$
  - ▶ Emits observation  $o_{t+1}$
  - ▶ Emits scalar reward  $r_{t+1}$

# State



- ▶ Experience is a sequence of observations, actions, rewards

$$o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t$$

- ▶ The **state** is a summary of experience

$$s_t = f(o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t)$$

- ▶ In a fully observed environment

$$s_t = f(o_t)$$

# Major Components of an RL Agent

- ▶ An RL agent may include one or more of these components:
  - ▶ **Policy**: agent's behaviour function
  - ▶ **Value function**: how good is each state and/or action
  - ▶ **Model**: agent's representation of the environment



# Policy

- ▶ A **policy** is the agent's behaviour
- ▶ It is a map from state to action:
  - ▶ Deterministic policy:  $a = \pi(s)$
  - ▶ Stochastic policy:  $\pi(a|s) = \mathbb{P}[a|s]$

# Value Function

- ▶ A **value function** is a prediction of future reward
  - ▶ “How much reward will I get from action  $a$  in state  $s$ ?”
- ▶  **$Q$ -value function** gives expected total reward
  - ▶ from state  $s$  and action  $a$
  - ▶ under policy  $\pi$
  - ▶ with discount factor  $\gamma$

$$Q^\pi(s, a) = \mathbb{E} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

# Value Function

- ▶ A **value function** is a prediction of future reward
  - ▶ “How much reward will I get from action  $a$  in state  $s$ ?”
- ▶  **$Q$ -value function** gives expected total reward
  - ▶ from state  $s$  and action  $a$
  - ▶ under policy  $\pi$
  - ▶ with discount factor  $\gamma$

$$Q^\pi(s, a) = \mathbb{E} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

- ▶ Value functions decompose into a Bellman equation

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

# Optimal Value Functions

- ▶ An optimal value function is the maximum achievable value

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

# Optimal Value Functions

- ▶ An optimal value function is the maximum achievable value

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

- ▶ Once we have  $Q^*$  we can act optimally,

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

# Optimal Value Functions

- ▶ An optimal value function is the maximum achievable value

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

- ▶ Once we have  $Q^*$  we can act optimally,

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

- ▶ Optimal value maximises over all decisions. Informally:

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

# Optimal Value Functions

- ▶ An optimal value function is the maximum achievable value

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

- ▶ Once we have  $Q^*$  we can act optimally,

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

- ▶ Optimal value maximises over all decisions. Informally:

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

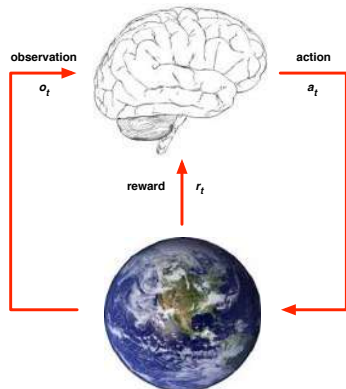
- ▶ Formally, optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

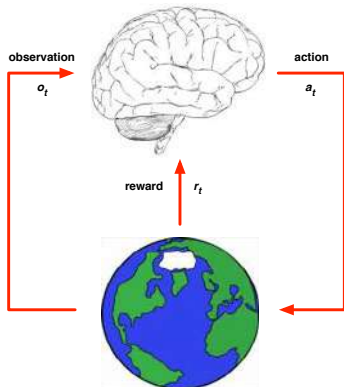
# Value Function Demo



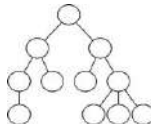
# Model



# Model



- ▶ **Model** is learnt from experience
- ▶ Acts as proxy for environment
- ▶ Planner interacts with model
- ▶ e.g. using lookahead search



# Approaches To Reinforcement Learning

## Value-based RL

- ▶ Estimate the **optimal value function**  $Q^*(s, a)$
- ▶ This is the maximum value achievable under any policy

## Policy-based RL

- ▶ Search directly for the **optimal policy**  $\pi^*$
- ▶ This is the policy achieving maximum future reward

## Model-based RL

- ▶ Build a model of the environment
- ▶ Plan (e.g. by lookahead) using model

# Deep Reinforcement Learning

- ▶ Use deep neural networks to represent
  - ▶ Value function
  - ▶ Policy
  - ▶ Model
- ▶ Optimise loss function by stochastic gradient descent

# Outline

Introduction to Deep Learning

Introduction to Reinforcement Learning

**Value-Based Deep RL**

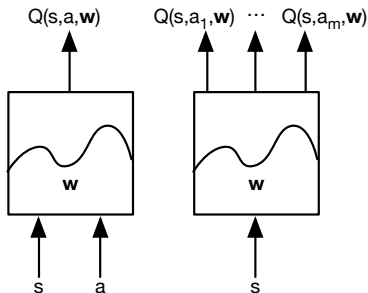
Policy-Based Deep RL

Model-Based Deep RL

# Q-Networks

Represent value function by **Q-network** with weights  $\mathbf{w}$

$$Q(s, a, \mathbf{w}) \approx Q^*(s, a)$$



# Q-Learning

- ▶ Optimal Q-values should obey Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

- ▶ Treat right-hand side  $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$  as a target
- ▶ Minimise MSE loss by stochastic gradient descent

$$l = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

# Q-Learning

- ▶ Optimal Q-values should obey Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

- ▶ Treat right-hand side  $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$  as a target
- ▶ Minimise MSE loss by stochastic gradient descent

$$l = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

- ▶ Converges to  $Q^*$  using table lookup representation



# Q-Learning

- ▶ Optimal Q-values should obey Bellman equation

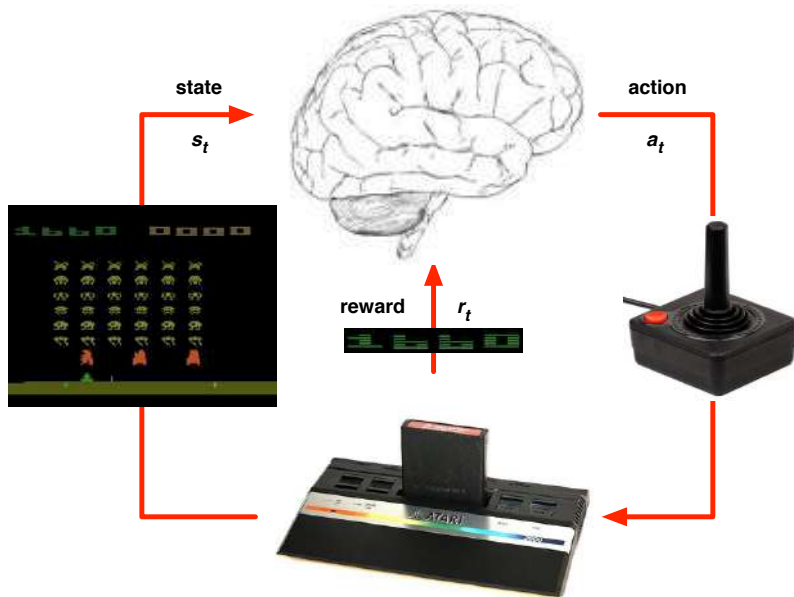
$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

- ▶ Treat right-hand side  $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$  as a target
- ▶ Minimise MSE loss by stochastic gradient descent

$$l = \left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

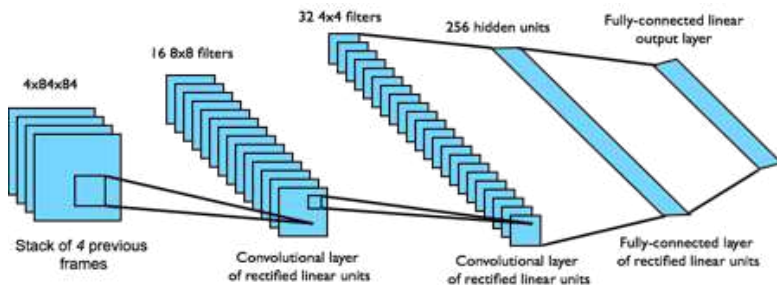
- ▶ Converges to  $Q^*$  using table lookup representation
- ▶ But **diverges** using neural networks due to:
  - ▶ Correlations between samples
  - ▶ Non-stationary targets

# Deep Reinforcement Learning in Atari



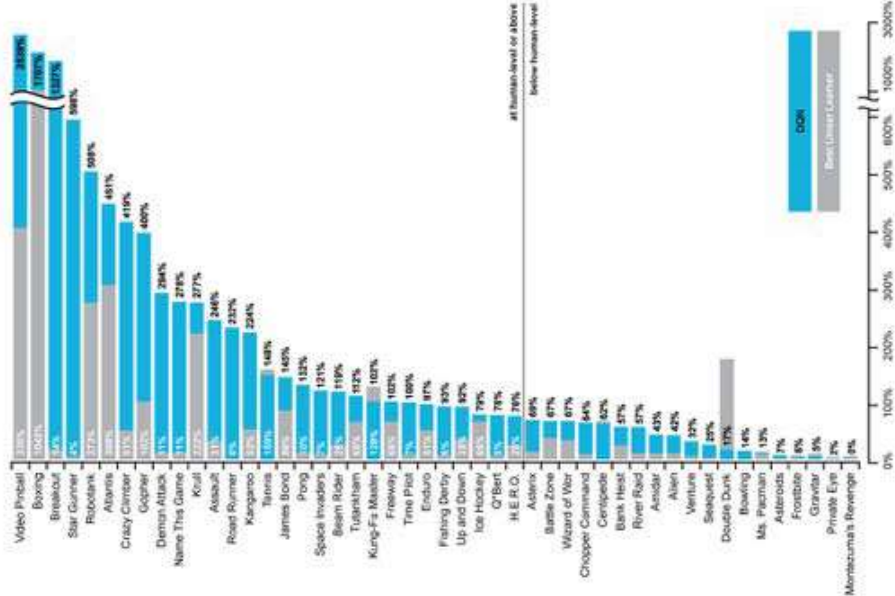
# DQN in Atari

- ▶ End-to-end learning of values  $Q(s, a)$  from pixels  $s$
- ▶ Input state  $s$  is stack of raw pixels from last 4 frames
- ▶ Output is  $Q(s, a)$  for 18 joystick/button positions
- ▶ Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

# DQN Results in Atari



# DQN Atari Demo

DQN paper

[www.nature.com/articles/nature14236](http://www.nature.com/articles/nature14236)

DQN source code:

[sites.google.com/a/deepmind.com/dqn/](https://sites.google.com/a/deepmind.com/dqn/)



# Conclusion

- ▶ General, stable and scalable RL is now possible
- ▶ Using deep networks to represent value, policy, model
- ▶ Successful in Atari, Labyrinth, Physics, Poker, Go
- ▶ Using a variety of deep RL paradigms