RL: studying (sequential decision-making)

→ sequential decision making
AI: → learning
→ deep computational graph
RL: study sequential decision-making
(general framework)

The Deep RL Boom

TD-Gammon ( Tesauro, 1989-1995)

↓

Slot car driving ( Lange & Riedmiller, 2012)

↓

Arcade Learning Environment (Bellemare 2013)

↓

Deep Q-Network (2013, 2015)

↓

Trust region policy optimization ( Schulman, 2015)

↓

End-to-end training ( Levine, 2015)

On real robots

↓

<u>AlphaGo</u>   ( 2015 )

---

## RL + DNN ?

① Sparse / delayed feedback

② data distribution is non-stationary

— determined by the agent's actions

- Most of the DNN theory no longer apply
- Exploration vs exploitation
- Local minima

③ Training DNN with RL was thought to be

inherently unstable.

( Tsitsikis & Van Roy (1997)

# DQN (Deep Q Networks)

CNN $\xrightarrow{\text{represent}}$ Action-value function (Q)

Tabular Q-Learning:

① Start with a guess for each $Q(s, a)$

② interact with environment using policy based on Q

③ Updates (based on Bellman Equations)

$(S_t, a_t, r_t, S_{t+1})$ experience

$\varepsilon$-greedy

$$Q(S_t, a_t) \longleftarrow Q(S_t, a_t) + \alpha_t \left( \underbrace{\left( r_t + \gamma \max_a Q(S_{t+1}, a) \right)}_{\text{target}} - Q(S_t, a_t) \right)$$

Problems: (DNN/CNN $\longrightarrow$ Q)

① Correlation between successive updates

② Correlation between $Q(S_t, a_t)$ and the target

High-level idea: Q-Learning $\xrightarrow[\text{Like}]{\text{Look}}$ Supervised learning

(RL)

Apply Q-update on batches of past experience

instead of online

1. Experience replay ( Lin, 1993)

2. Previously used for better data efficiency

3. Makes the data distribution more stationary

Use an old set of weights to compute the

target ( target network)

Keeps the target function from changing too quickly

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left( r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s,a; \theta_i) \right)^2$$

target

# Neural Fitted Q Iteration

## NFQ (Riedmiller, 2005)

## Target Network Intuition

- Changing the value of one action will change the value of other actions and similar states.
- The network can end up chasing its own tail because of bootstrapping.
- Somewhat surprising fact - bigger networks are less prone to this because they alias less.

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left( \underbrace{r + \gamma \max_{a'} Q(s', a'; \theta_i^-)}_{\text{target}} - Q(s, a; \theta_i) \right)^2$$

$s$ $\qquad$ $s'$

DeepMind

because they can use action repeats but these are two successive

MacBook Pro

# DQN Training Algorithm

**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 $a_t = \mathrm{argmax}_a Q(\phi(s_t), a; \theta)$

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

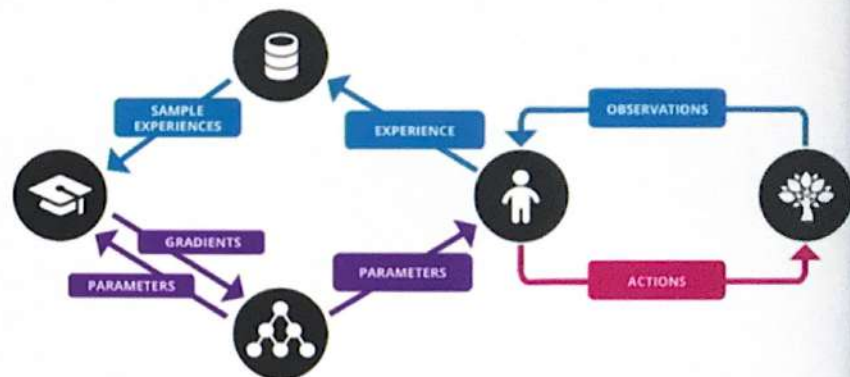        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$

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

        Perform a gradient descent step on $\left( y_j - Q(\phi_j, a_j; \theta) \right)^2$ with respect to the network parameters $\theta$

        Every $C$ steps reset $\hat{Q} = Q$

    **End For**

**End For**

SAMPLE EXPERIENCES — EXPERIENCE — OBSERVATIONS — GRADIENTS — PARAMETERS — PARAMETERS — ACTIONS

DeepMind

MacBook Pro

# Neural Fitted Q iteration

## Neural Fitted Q Iteration

- NFQ (Riedmiller, 2005) trains neural networks with Q-learning.
- Alternates between collecting new data and fitting a new Q-function to all previous experience with batch gradient descent.

```
NFQ_main() {
input: a set of transition samples D; output: Q-value function Q_N
    k=0
    init_MLP() → Q_0;
    Do {
        generate_pattern_set P = {(input^l, target^l), l = 1, ..., #D} where:
            input^l = s^l, u^l,
            target^l = c(s^l, u^l, s'^l) + γ min_b Q_k(s'^l, b)
        Rprop_training(P) → Q_{k+1}
        k:= k+1
    } WHILE (k < N)
```

- DQN can be seen as an online variant of NFQ.

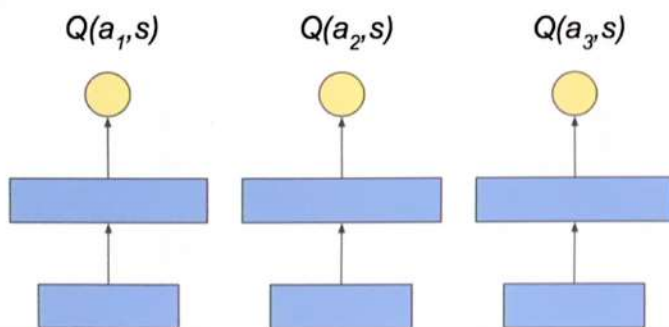DeepMind

similar method to dqn was neural fitted cue iteration so
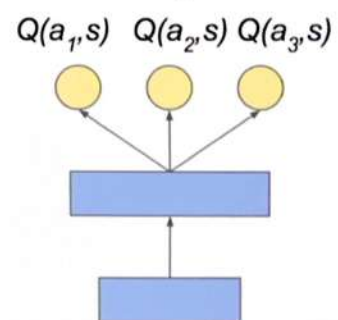
MacBook Pro

Lin's Networks

# Lin's Networks

- Long-Ji Lin's thesis "Reinforcement Learning for Robots using Neural Networks" (1993) also trained neural nets with Q-learning.
- Introduced experience replay among other things.
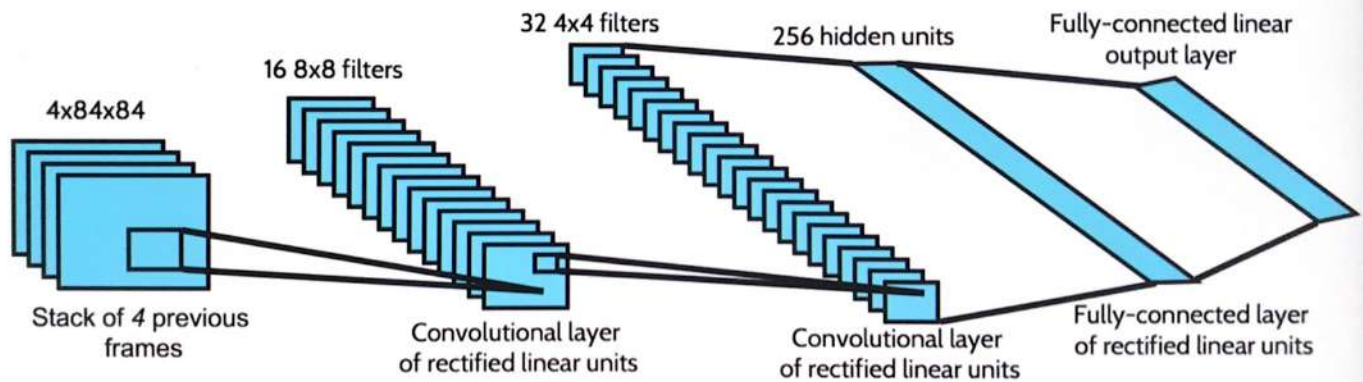- Lin's networks did not share parameters among actions.

**Lin's architecture**

$Q(a_1,s)$     $Q(a_2,s)$     $Q(a_3,s)$

**DQN**

$Q(a_1,s)$  $Q(a_2,s)$  $Q(a_3,s)$

DeepMind

# ATARI Network Architecture

- Convolutional neural network architecture:
  - History of frames as input.
  - One output per action - expected reward for that action $Q(s, a)$.
  - Final results used a slightly bigger network (3 convolutional + 1 fully-connected hidden layers).



**4x84x84**

**16 8x8 filters**

**32 4x4 filters**

**256 hidden units**

**Fully-connected linear output layer**

Stack of 4 previous frames

Convolutional layer of rectified linear units

Convolutional layer of rectified linear units

Fully-connected layer of rectified linear units

one output per action so not only by today's standards

MacBook Pro

# DQN/ Mini-batch Q-learning

**(Pros)**:

- Robust

- GPU friendly

**(Cons)**:

- Slow

- less RNN friendly ( less successful
  in 3D environment)

## Beyond DQN

① fast training

② on or off-policy methods

③ flexibility ⟵ discrete or continuous actions
            feedforward
            recurrent models

④ Asynchronous Methods for DRL

## AsyncRL

Parallel actor-learners (CPU threads)

online asynchronous updates

(Recht 2011, Lian 2015)

RL algorithm:
- on-policy / off-policy
- value-based / policy-based

## 1-step Q-learning

Parallel actor-learners compute online 1-step update.

$$y \longleftarrow r + \gamma \max_{a'} Q(s', a'; \theta^-)$$

$$\Delta\theta \longleftarrow \Delta\theta + \frac{\partial (y - Q(s, a; \theta))^2}{\partial \theta}$$

# N-step Q-learning

Q-learning with a uniform mixture of backups
of length 1 through N.

$$r_t, \; r_{t+1}, \; r_{t+2}, \; \ldots \; , \; r_{t+n}. \quad \text{Max} \; \chi_a \; Q(a, S_{t+N+1})$$

$$y \leftarrow \sum_{k=0}^{N-1} \gamma^k \, r_{t+k} + \gamma^N \max_{a'} Q(S_{t+N}, a'; \theta^-)$$

$$\Delta\theta \leftarrow \Delta\theta + \frac{\partial (y - Q(S_t, a_t; \theta))^2}{\partial \theta}$$

Variation of "Incremental multi-step Q-learning"

(Peng & Williams 1995)

# Asynchronous Advantage Actor-critic
## (A3C)

① Agent Learns
- a policy
- a state value function

② bootstrapped n-step returns

$\longrightarrow$ reduce variance

over REINFORCE with a baseline

③ Policy gradient multiplied by an estimate of the advantage

$$\nabla_\theta \log \pi(a_t \mid s_t, \theta) \left( \sum_{k=0}^{N} \gamma^k r_{t+k} + \gamma^{N+1} V(S_{t+N+1}) - V(S_t) \right)$$

④ Critic/value function is trained with n-step TD learning (by minimizing the MSE)

$$\sum_{k=0}^{N} \gamma^k r_{t+k} + \gamma^{N+1} V(S_{t+N+1} ; \theta^-) - V(S_t ; \theta))^2$$

> " A3c tends to dominate the value-based Methods "

## A3C - ATARI Results

| Method | Training Time | Mean | Median |
|---|---|---|---|
| DQN | 8 days on GPU | 121.9% | 47.5% |
| Gorilla | 4 days, 100 machines | 215.2% | 71.3% |
| D-DQN | 8 days on GPU | 332.9% | 110.9% |
| Dueling D-DQN | 8 days on GPU | 343.8% | 117.1% |
| Prioritized DQN | 8 days on GPU | 463.6% | 127.6% |
| **A3C, FF** | 1 day on CPU | 344.1% | 68.2% |
| **A3C, FF** | 4 days on CPU | 496.8% | 116.6% |
| **A3C, LSTM** | 4 days on CPU | 623.0% | 112.6% |

DeepMind

Atari and if you train for four days then you

MacBook Pro

Pros of N-step Methods:

① faster reward propagation

② No need for target network

③ easier training of RNN

A3C —— Pros ( fast , RNN friendly)
                ( stable, scalable)

              Cons ( Not GPU friendly)


IMPALA :   distributed  deep  RL ,(2018)


Scale  up  (A3C)

" V - Trace " algorithm

" DM Lab - 30  Task Set "
              (Multi-)


DQN  is  more  stable  then  A3C


off-policy  harms  actor - critic  method

# Deep RL

① Minibatch training on GPUs

② Deep ResNet and LSTMs

③ Adam / RMSProp optimizers

Deep RL

# Practical Advice - Getting Started

- Start with a simple problem.
    - Something solvable in under a minute on your local machine.
    - Make it similar to the problem you really want to solve.
    - Ideally it should have knobs for controlling its difficulty.
- Plot the training curves (averaged over multiple episodes).
- Visualize the policy.
- Visualize the value function.
- Visualize everything you can think of.

DeepMind

has some of the properties of the real
problem you want to

# Practical Advice - Neural Nets

- Doing early experiments with a small network can help iterate faster.
  - This can also backfire (DQN and target networks).
- Reasonable strategy:
  - Run a few progressively larger nets to find what's sufficient for experimenting.
  - Periodically try larger nets to max out performance and verify assumptions.
- Be careful with initialization:
  - Visualize the initial policy to make sure it gets some rewards.
- Try RMSProp and/or Adam.
- Test deep learning tricks before incorporating them: dropout, batch norm, etc.
- See John Schulman's excellent guide - http://joschu.net/docs/nuts-and-bolts.pdf

DeepMind