

# Deep Reinforcement Learning: A brief introduction

A/Prof Richard Yi Da Xu

Yida.Xu@uts.edu.au

Wechat: aubedata

<https://github.com/roboticcam/machine-learning-notes>

University of Technology Sydney (UTS)

February 18, 2018

- ▶ A video from Google DeepMind's Deep Q-learning playing Atari Breakout:  
<https://www.youtube.com/watch?v=TmPfTpjtdgg>
- ▶ Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." *arXiv preprint arXiv:1312.5602* (2013).
- ▶ code is also available  
<https://github.com/kuz/DeepMind-Atari-Deep-Q-Learner>

## N.B.

- ▶ Apologies for those have seen it before

**significance** of this demo shows it's possible to use Neural Network to learn how to play a game, based on:

- ▶ sequences of screen images
- ▶ scores the game receives
- ▶ goal is to learn the best policy for **actions** to take

Surely you **don't need a menu** to learn how to play Atari. i.e., it's **model-free**!

Forget about the Neural network for a second, how is Reinforcement Learning (RL) different to conventional supervised learning?

- ▶ No data label like supervised learning, i.e., no “*best-action-for-that-screen*” label
- ▶ only **reward** signal
- ▶ feedback in **delayed**, not instantaneous
- ▶ data are not i.i.d., (consecutive frames are similar)
- ▶ agent's actions affects the subsequent data it receives.

Let's get started with some RL background.

another way to look at it:

- ▶ RL uses training information that **evaluates the actions** taken rather than **instructs by giving correct actions**.
- ▶ a need for active exploration: explicit trial-and-error search for good behavior.
- ▶ **purely evaluative feedback** indicates how good the action taken is, but not whether it is the best or the worst action possible.
- ▶ **purely instructive feedback** indicates correct action to take, independently of the action actually taken. [supervised learning](#)

- ▶ **marketing**

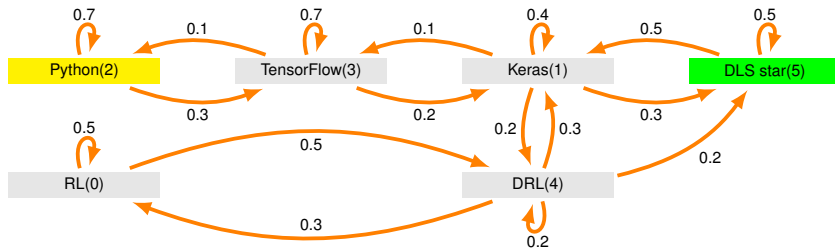
customer's attributes  $s$ , marketing actions  $a$ , customer signs up  $r$

- ▶ **drone control**

all available sensor data  $a$ , controls  $s$ , not crashing  $r$

- ▶ **chatbot**

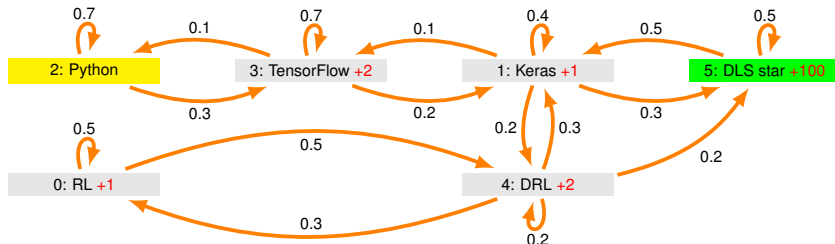
conversations to-date  $s$ , things that a robot will say  $a$ , customer satisfaction  $r$



- ▶ one may start from **python** and generate sequences with transition probabilities to end up in **DLS star**. examples:
  - ▶ Python, Python, Python, TensorFlow, Keras, DLS star
  - ▶ Python, Python, Python, TensorFlow, TensorFlow, Keras, DRL, DRL RL, DLS star
  - ▶ Python, Python, TensorFlow, TensorFlow, Keras, DRL, DLS star
  - ▶ The question is: how we may able to measure “how good” each path? ...

# Markov Reward Process

Let's add some rewards to being at each of the state:



What we care is the **total return**  $G_t$ : sum of **discounted** reward from time-step  $t$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad \text{where } \gamma \in [0, 1]$$

note that  $G_t$  is a random variable

**exercise** what happens when  $\gamma = 0$  and  $\gamma = 1$

# Markov Random Process: State value function

- ▶ state value function  $v(s)$
- ▶ is expected total return starting from state  $s$

$$\begin{aligned}v(s) &\equiv v(S_t) = \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [G_t | s_t = s] \\&= \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]\end{aligned}$$

- ▶  $v(s)$  indicates the long future value, if you are currently at state  $s$



# Markov Random Process: Bellman Equation

- ▶ state value function  $v(s)$  is expected total return starting from state  $s$

$$\begin{aligned}v(s) &= \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [G_t | s_t = s] \\&= \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s] \\&= \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [R_{t+1} + \gamma \underbrace{(R_{t+2} + \gamma R_{t+3} + \dots)}_{G_{t+1}} | s_t = s] \\&= \mathbb{E}_{s_{t+1}, s_{t+2}, \dots} [R_{t+1} + \gamma G_{t+1} | s_t = s]\end{aligned}$$

- ▶ we know that,  $\mathbb{E}_{X,Y}[r(X) + r(X, Y)] = \mathbb{E}_X[r(X) + \mathbb{E}_Y[r(X, Y)]]$

$$\begin{aligned}v(s) &= \mathbb{E}_{s_{t+1}} [R_{t+1} + \gamma \mathbb{E}_{s_2, s_3, \dots} [G_{t+1} | s_{t+1}] | s_t = s] \\&= \mathbb{E}_{s_{t+1}} [R_{t+1} + \gamma v(s_{t+1}) | s_t = s]\end{aligned}$$

- ▶ **Bellman equations:** value of the current state,  $v(s)$  breaks up into (1) **immediate** and (2) **future** rewards.

# Bellman Equation in matrix form

$$v(s_t) \equiv v(s) = \mathbb{E}_{s_1} [R_{t+1} + \gamma v(S_{t+1}) | s_t = s]$$

► say  $s \in \{1, \dots, n\}$ :

$$\underbrace{v(s_t = 1)}_{v(1)} = \mathbb{E}_{s_1} [\underbrace{R_{t+1}(s_t = 1)}_{R_1} + \gamma v(S_{t+1}) | s_t = 1]$$

$$v(s_t = 2) = \mathbb{E}_{s_1} [R_{t+1}(s_t = 2) + \gamma v(S_{t+1}) | s_t = 2]$$

...

take the first line,

$$v(1) = \mathbb{E}_{s_1} [R_1 + \gamma v(S_{t+1}) | s_t = 1]$$

$$= R_1 + \gamma \mathbb{E} [v(S_{t+1}) | s_t = 1]$$

$$= R_1 + \gamma \left( \sum_{s_{t+1}=1}^n v(s_{t+1}) \Pr(1 \rightarrow s_{t+1}) \right)$$

$$= R_1 + \gamma \left( \sum_{j=1}^n v(j) \Pr(1 \rightarrow j) \right)$$

...

$$\Rightarrow v(k) = R_k + \gamma \left( \sum_{j=1}^n v(j) \Pr(k \rightarrow j) \right)$$

## Bellman Equation in matrix form (2)

$$\begin{aligned}v(k) &= R_k + \gamma \left( \sum_{j=1}^n v(j) \Pr(k \rightarrow j) \right) \\&= R_k + \gamma \mathcal{P}_{k,:}^\top \mathbf{v} \\ \Rightarrow \mathbf{v} &= \mathbf{R} + \gamma \mathcal{P} \mathbf{v} \\ \Rightarrow \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} &= \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{1,1} & \dots & \mathcal{P}_{1,n} \\ \vdots & & \vdots \\ \mathcal{P}_{n,1} & \dots & \mathcal{P}_{n,n} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}\end{aligned}$$

the solution to **MRP** is straight forward:

$$\begin{aligned}\mathbf{v} &= \mathbf{R} + \gamma \mathcal{P} \mathbf{v} \\ (I - \gamma \mathcal{P}) \mathbf{v} &= \mathbf{R} \\ \mathbf{v} &= (I - \gamma \mathcal{P})^{-1} \mathbf{R}\end{aligned}$$

# Markov Decision Process (MDP)

- ▶ now agent has **actions**
- ▶ concept of **policy**  $\pi$ : take a state  $s_t$  as input and decides an action  $a_t$

$$\pi(a|s) = \Pr(A_t = a | S_t = s)$$

- ▶ a policy is time-invariant (or stationary) and stochastic
- ▶ next state for an agent, now also depends on its action taken:

$$\mathcal{P}_{s \rightarrow s'}^a = \Pr(S_{t+1} = s' | S_t = s, A_t = a)$$

- ▶ multiple transition matrix  $\mathcal{P}$  each depends on the  $a$  taken
- ▶ once fixed  $\pi$ , MDP becomes MRP with transition probability  $\mathcal{P}_{s \rightarrow s'}^\pi$ :

$$\mathcal{P}_{s \rightarrow s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{P}_{s \rightarrow s'}^a$$

- ▶ is the expected return starting from state  $s$ , and then follow policy  $\pi$ :

$$v^\pi(s) = \mathbb{E}[G_t | s_t = s]$$

- ▶ for example,  $v^\pi(s_0)$ :

$$v^\pi(s_0) = \mathbb{E}_{s_1, s_2, s_3, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t \underbrace{r_t(s_t, \pi(s_t), s_{t+1})}_{\text{reward}} | s_0 \right]$$

- ▶ the **optimal state value** is computed using **best policy**:

$$V^*(s_0) = \max_{\pi} \left( \mathbb{E}_{s_1, s_2, s_3, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) | s_0 \right] \right)$$

$$\begin{aligned}
 V^*(s_0) &= \max_{\pi} \left( \mathbb{E}_{s_1, s_2, s_3, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t \underbrace{r_t(s_t, \pi(s_t), s_{t+1})}_{\text{red}} \middle| s_0 \right] \right) \\
 &= \max_{\pi} \left( \mathbb{E}_{s_1} \left[ \mathbb{E}_{s_2, s_3, \dots} \left[ \gamma^0 r_0(s_0, \pi(s_0), s_1) + \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \right] \middle| s_0 \right] \right) \\
 &= \max_{\pi} \left( \mathbb{E}_{s_1} \left[ \gamma^0 r_0(s_0, \underbrace{\pi(s_0)}_{a_0}, s_1) + \mathbb{E}_{s_2, s_3, \dots} \left[ \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right] \middle| s_0 \right] \right)
 \end{aligned}$$

**since**  $a_t = \pi(s_t) \implies \max_{\pi} f(\pi(s_t)) = \max_{a_t} f(\pi(s_t))$  **and**  $\pi$  still takes care of rest of  $a_1, \dots$

$$= \max_{a_0} \left( \mathbb{E}_{s_1} \left[ \gamma^0 r_0(s_0, \pi(s_0), s_1) + \max_{\pi} \left( \mathbb{E}_{s_2, s_3, \dots} \left[ \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right] \right) \middle| s_0 \right] \right)$$

we like the first  $\gamma$  to be 1

$$= \max_{a_0} \left( \mathbb{E}_{s_1} \left[ \underbrace{\gamma^0}_{=1} r_0(s_0, \pi(s_0), s_1) + \underbrace{\gamma \max_{\pi} \left( \mathbb{E}_{s_2, s_3, \dots} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right] \right)}_{V^*(s_1)} \middle| s_0 \right] \right)$$

$$\begin{aligned} V^*(s_0) &= \max_{\pi} \left( \mathbb{E}_{s_1, s_2, s_3, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_0 \right] \right) \\ &= \max_{a_0} \left( \mathbb{E}_{s_1} \left[ r_0(s_0, \pi(s_0), s_1) + \underbrace{\gamma \max_{\pi} \left( \mathbb{E}_{s_2, s_3, \dots} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right] \right)}_{V^*(s_1)} \middle| s_0 \right] \right) \\ &= \max_{a_0} \left( \mathbb{E}_{s_1} \left[ r_0(s_0, \pi(s_0), s_1) + \gamma V^*(s_1) \middle| s_0 \right] \right) \end{aligned}$$

# Bellman's equation

$$V^*(s) = \max_a \left( \mathbb{E}_{s'} \left[ r(s, \pi(s), s') + \gamma V^*(s') \right] \middle| s \right)$$

- ▶ drop  $|s$ , one has:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{s'} \left[ r(s, \pi(s), s') + \gamma V^\pi(s') \right] \\ \implies V^\pi(s) + \eta V^\pi(s) &= V^\pi(s) + \eta \left( \mathbb{E}_{s'} \left[ r(s, \pi(s), s') + \gamma V^\pi(s') \right] \right) \\ \implies V^\pi(s) &= V^\pi(s) + \eta \left( \mathbb{E}_{s'} \left[ r(s, \pi(s), s') + \gamma V^\pi(s') \right] - V^\pi(s) \right) \end{aligned}$$

- ▶ instead of compute this expectation, in **each iteration**  $t$ , we sample a new state  $\tilde{s}' \sim \Pr(s' | \dots)$

$$V_{t+1}^\pi(s) = V_t^\pi(s) + \eta \left( r(s, \pi(s), \tilde{s}') + \gamma V_t^\pi(\tilde{s}') - V_t^\pi(s) \right)$$

- ▶ note that the last equation is called **temporal difference**



# Bellman's equation: Three ways of solving it

$$V^\pi(s) = \mathbb{E}_\pi[G_t | s_t = s]$$

— could be approximated by **Monte-carlo**, i.e., sample  $s_{t+1}, s_{t+2}, \dots$  and compute  $G_t$

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

— could be approximated by **Temporal Difference**

$$= \sum_a \pi(a|s) \sum_{s'} \mathcal{P}_{s \rightarrow s'}^a \left[ r(s, \pi(s), s') + \gamma V^\pi(s') \right]$$

— could be solved exactly by **Dynamic programming**

# Action-value (Q) function

- ▶ action-valued function  $Q^\pi(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a, \pi]$ :
- ▶ expected total return starting from state  $s$ , taking action  $a$ , and then follow policy  $\pi$
- ▶ Stochastic policy  $\pi$ :

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)}[Q^\pi(s, a)]$$

- ▶ deterministic policy:

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

- ▶ from before;

$$\begin{aligned} V^*(s) &= \max_a \left( \mathbb{E}_{s'} \left[ r(s, a, s') + \gamma \underbrace{V^*(s')} \middle| s \right] \right) \\ &= \max_a \left( \underbrace{\mathbb{E}_{s'} \left[ r(s, a, s') + \gamma \left( \max_{a'} Q^*(s', a') \right) \right] \middle| s}_{Q^*(s', a') \text{ by definition}} \right) \end{aligned}$$

- ▶ therefore:

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

# Action-value (Q) function

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

- drop  $|s, a$ , let's solve this by **temporal difference**:

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

$$\implies \textcolor{red}{Q}^\pi(\textcolor{red}{s}, \textcolor{red}{a}) + \eta Q^\pi(s, a) = \textcolor{red}{Q}^\pi(\textcolor{red}{s}, \textcolor{red}{a}) + \eta \left( \mathbb{E}_{s'} \left[ r(s, a, s') + \gamma \left( \max_{a'} Q^\pi(s', a') \right) \right] \right)$$

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

- instead of compute this expectation, in **each iteration**  $t$ , we sample a new state  $(\tilde{s}', \tilde{a}) \sim \Pr(s', a | \dots)$ .

**Q-Learning**: recursively:

$$Q(s, \tilde{a}) = Q(s, \tilde{a}) + \eta \left( \underbrace{r(s, \tilde{a}, \tilde{s}') + \gamma \left( \max_{a'} Q(\tilde{s}', a') \right)}_y - Q(s, \tilde{a}) \right)$$

- let  $\eta = 1$ :

$$Q(s, \tilde{a}) = r(s, \tilde{a}, \tilde{s}') + \gamma \left( \max_{a'} Q(\tilde{s}', a') \right)$$

- ▶ advantage function:

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

$$\begin{aligned} & \mathbb{E}_{a \sim \pi(s)}[A^\pi(s, a)] \\ &= \mathbb{E}_{a \sim \pi(s)}[Q^\pi(s, a) - V^\pi(s)] \\ &= \mathbb{E}_{a \sim \pi(s)}[Q^\pi(s, a)] - V^\pi(s) = 0 \end{aligned}$$

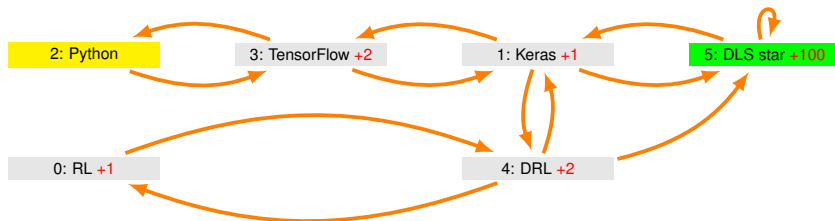
- ▶  $V$  measures how good in a particular state  $s$ .
- ▶  $Q$  measures value of choosing a particular action when in this state.
- ▶  $A$  measures **relative importance** of each action.

**Require:** choice of  $\gamma$                       Rewards matrix  $R$

```
1:  $Q \leftarrow \mathbf{0}$ 
2: for each episode do
3:   randomise initiate state  $s_0$ 
4:   while goal state not reached do
5:     select  $(a, s') \sim \text{Pr}(a, s' | \cdot)$ 
6:     compute  $\max_{a'} Q(s', a')$ 
7:      $Q(s, a) \leftarrow r(s, a, s') + \gamma (\max_{a'} Q(s', a'))$ 
8:      $s_t \leftarrow s_{t+1}$ 
9:   end while
10: end for
```

# Q-Learning example

We took the example from the Markov Reward Process example earlier:



- ▶ there is small immediate rewards by going from one module to another
- ▶ you get a final large reward by becoming DLS star
- ▶ let  $\gamma = 0.5$
- ▶ in this special example,  $a = s'$ , i.e., the action is to turn into the next state (module of studies).
- ▶ assume equal probabilities for all edges.

# Q-Learning example: episode 1

$$R = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & - & - & - & - & 2 & - \\ \text{Ke}(1) & - & - & - & 2 & 2 & 100 \\ \text{Py}(2) & - & - & - & 2 & - & - \\ \text{TF}(3) & - & 1 & 0 & - & - & - \\ \text{DRL}(4) & 1 & 1 & - & - & - & 100 \\ \text{DLS}^*(5) & - & 1 & - & - & - & 100 \end{matrix}$$

before

$$Q = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{Ke}(1) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{Py}(2) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{TF}(3) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DRL}(4) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DLS}^*(5) & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

after

$$Q = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{Ke}(1) & 0 & 0 & 0 & 0 & 0 & 100 \\ \text{Py}(2) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{TF}(3) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DRL}(4) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DLS}^*(5) & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

- ▶  $s \sim \text{Pr}(s|\cdot) = 1$ , i.e. Keras
- ▶ at  $s = 1$ , it has **allowable actions**: go to state  $\{3, 4, 5\}$ , i.e.,  $a \in \{3, 4, 5\}$
- ▶  $(a, s') \sim \text{Pr}(a, s'|\cdot) = (5, 5)$
- ▶ at  $s' = 5$ , it has **allowable actions**:  $a' \in \{1, 5\}$ :

$$\begin{aligned} Q(s, a) &= r(s, a, s') + \gamma \left( \max_{a'} Q(s', a') \right) \\ &= R(1, s' = 5) + 0.5 \max[Q(s' = 5, 1), Q(s' = 5, 5)] \\ &= 100 + 0.5 \times 0 = 100 \end{aligned}$$

- ▶ set  $s \leftarrow s' \implies s = 5$ , i.e., goal state, end

# Q-Learning example: episode 2, Iteration 1

$$R = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & - & - & - & - & 2 & - \\ \text{Ke}(1) & - & - & - & 2 & 2 & 100 \\ \text{Py}(2) & - & - & - & 2 & - & - \\ \text{TF}(3) & - & 1 & 0 & - & - & - \\ \text{DRL}(4) & 1 & 1 & - & - & - & 100 \\ \text{DLS}^*(5) & - & 1 & - & - & - & 100 \end{matrix}$$

before

$$Q = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{Ke}(1) & 0 & 0 & 0 & 0 & 0 & 100 \\ \text{Py}(2) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{TF}(3) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DRL}(4) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DLS}^*(5) & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

after

$$Q = \begin{matrix} S \downarrow, A \rightarrow & \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \text{RL}(0) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{Ke}(1) & 0 & 0 & 0 & 0 & 0 & 100 \\ \text{Py}(2) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{TF}(3) & 0 & 51 & 0 & 0 & 0 & 0 \\ \text{DRL}(4) & 0 & 0 & 0 & 0 & 0 & 0 \\ \text{DLS}^*(5) & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

- ▶  $s \sim \Pr(s|\cdot) = 3$
- ▶ at  $s = 3$ , it has **allowable actions**: go to state  $\{1, 2\}$ , i.e.,  $a \in \{1, 2\}$
- ▶  $(a, s') \sim \Pr(a, s'|\cdot) = (1, 1)$
- ▶ at  $s' = 1$ , it has **allowable actions**:  $a' \in \{3, 4, 5\}$ :

$$\begin{aligned} Q(s, a) &= r(s, a, s') + \gamma \left( \max_{a'} Q(s', a') \right) \\ &= R(3, 1) + 0.5 \max[Q(1, 3), Q(1, 4), Q(1, 5)] \\ &= 1 + 0.5 \times 100 = 51 \end{aligned}$$

- ▶ set  $s \leftarrow s' \implies s = 1$ , i.e., **not** a goal state, keep on going



# Q-Learning example: episode 2, Iteration 2

$$R = \begin{matrix} S \downarrow, A \rightarrow \\ \text{RL}(0) \\ \text{Ke}(1) \\ \text{Py}(2) \\ \text{TF}(3) \\ \text{DRL}(4) \\ \text{DLS}^*(5) \end{matrix} \begin{pmatrix} \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \begin{matrix} - & - & - & - & 2 & - \\ - & - & - & 2 & 2 & 100 \\ - & - & - & 2 & - & - \\ - & 1 & 0 & - & - & - \\ 1 & 1 & - & - & - & 100 \\ - & 1 & - & - & - & 100 \end{matrix} \end{pmatrix}$$

before

$$Q = \begin{matrix} S \downarrow, A \rightarrow \\ \text{RL}(0) \\ \text{Ke}(1) \\ \text{Py}(2) \\ \text{TF}(3) \\ \text{DRL}(4) \\ \text{DLS}^*(5) \end{matrix} \begin{pmatrix} \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 51 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \end{pmatrix}$$

after

$$Q = \begin{matrix} S \downarrow, A \rightarrow \\ \text{RL}(0) \\ \text{Ke}(1) \\ \text{Py}(2) \\ \text{TF}(3) \\ \text{DRL}(4) \\ \text{DLS}^*(5) \end{matrix} \begin{pmatrix} \text{RL}(0) & \text{Ke}(1) & \text{Py}(2) & \text{TF}(3) & \text{DRL}(4) & \text{DLS}^*(5) \\ \begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 51 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \end{pmatrix}$$

- ▶  $s = 1$  from previous iteration
- ▶ at  $s = 1$ , it has **allowable actions**: go to state  $\{3, 4, 5\}$ , i.e.,  $a \in \{3, 4, 5\}$
- ▶  $(a, s') \sim \Pr(a, s' | \cdot) = (5, 5)$
- ▶ at  $s' = 5$ , it has **allowable actions**:  $a' \in \{1, 5\}$ :

$$Q(s, a) = r(s, a, s') + \gamma \left( \max_{a'} Q(s', a') \right)$$

$$\begin{aligned} Q(1, 5) &= R(1, 5) + 0.5 \max[Q(5, 1), Q(5, 5)] \\ &= 100 + 0.5 \times 0 = 100 \end{aligned}$$

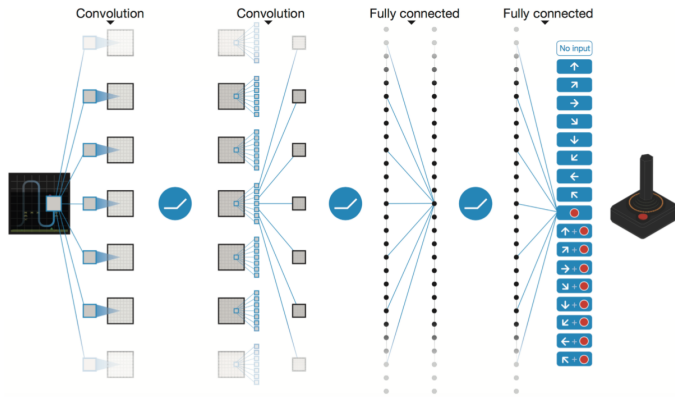
- ▶ set  $s \leftarrow s' \implies s = 5$ , i.e., goal state, end the state-action table gets updated until convergence.

# On the setting of Atari

- ▶ the states are far too many!
- ▶ need a **function approximator** to estimate the action-value function,  
 $Q(s, a|\theta) \approx Q^*(s, a)$
- ▶ guess what? Deep Neural Network helps!

# Represent $Q(s, a)$ using neural networks

- ▶ The figure below represents a row of the Q function table earlier:



Conv [16]  $\rightarrow$  ReLU  $\rightarrow$  Conv [32]  $\rightarrow$  ReLU  $\rightarrow$  FC [256]  $\rightarrow$  ReLU  $\rightarrow$  FC [|A|]

- ▶ these are **not** softmax functions.

# abstracted algorithm for Deep Q-Learning

**Require:** Initialize an empty replay memory

**Require:** Initialize the DQN weights  $\theta$

- 1: **for** each episode **do**
- 2:   **for**  $t = 1, \dots, T$  **do**
- 3:     with probability  $\epsilon$  select  $\tilde{a}$  random action
- 4:     otherwise, select:

$$\tilde{a} = \max_a (Q^*(s, a|\theta))$$

- 5:     perform  $\tilde{a}$  and receive rewards  $r_t$  and state  $s'$ .
- 6:     add tuple  $(s, \tilde{a}, r_t, s')$  into replay memory
- 7:     Sample a mini-batch of tuples  $(s_j, a_j, r_j, s'_j)$  from the replay memory
- 8:     and perform stochastic gradient descent on the DQN, based on the loss function:

$$\left( \underbrace{r_j + \gamma (\max_{a'} Q(s'_j, a'|\theta^-))}_{y_j} - Q(s_j, a_j|\theta) \right)^2$$

- 9:   **end for**
- 10: **end for**

innovation

- ▶ freeze parameters of target network  $Q(s'_j, a'|\theta^-)$  for fixed number of iterations
- ▶ while updating the online network  $Q(s; a; \theta_i)$  by gradient descent

- ▶ same values  $\theta$  both to select and to evaluate an action:

$$\begin{aligned}y_j &= r_j + \gamma(\max_{a'} Q(s'_j, a' | \theta)) \\ &= r_j + \gamma(Q(s'_j, \arg \max_a Q(s'_j, a, \theta) | \theta))\end{aligned}$$

- ▶ more likely to select overestimated values
- ▶ resulting in overoptimistic value estimates
- ▶ the solution is:

$$y_j = r_j + \gamma(\max_{a'} Q(s'_j, \arg \max_a Q(s'_j, a, \theta) | \theta'))$$

- ▶ still estimating value of policy according to current values defined by  $\theta$
- ▶ use second set of weights  $\theta'$  to **fairly** evaluate value of this policy

- ▶ CNN and RNN are two of the building blocks in Deep Learning
- ▶ People have been putting them into many existing machine learning frameworks, and have generated many interesting stuff
- ▶ but there is plenty still needs to be explored!

