# Lecture 14
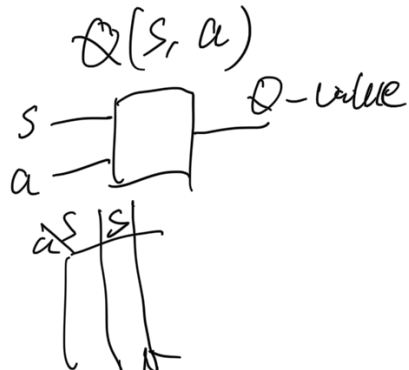
Keep Q-learning
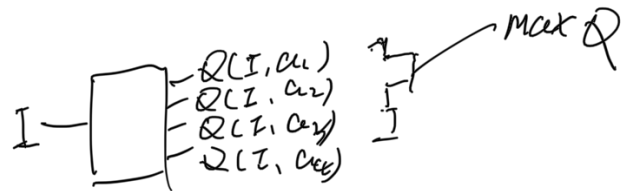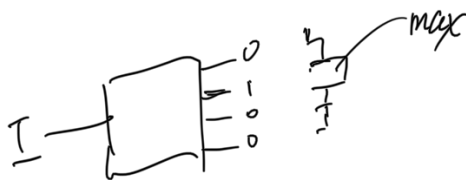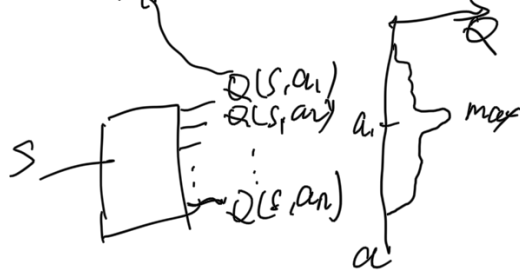
Look at inputs, outputs, Loss, dataset, structure

$Q(S, a)$

S ——[ ]—— Q-value
a ——

$a \in S \mid S$

when we use it, it is not very convenient

$S$ ——[ ]—— $Q_1$  . . . . $\begin{array}{c}S\\a_n\end{array}$——[ ]—— $Q_n$
$a_1$ ——

then get the max $Q_i$, that point to the action to take

$Q(S, a_1)$
$Q(S, a_2)$  $a_1$ —[ max ]— $Q$
$\vdots$
$Q(S, a_n)$
$a$

S ——[ ]——

$\begin{array}{c}0\\\vdots\\0\end{array}$ —[ max ]

I ——[ ]——

$Q(I, a_1)$
$Q(I, a_2)$ —[ max $Q$ ]
$Q(I, a_3)$
$Q(I, a_{44})$

I ——[ ]——

$\sqrt{ }$ random
$\textcircled{1}$  $(S, a, r, S')$

Episode, run from one random state, for several steps

record $\begin{cases} (S_1, a_1, r_1, S_2) & \text{step-index} \\ (S_2, a_2, r_2, S_3) \\ \quad \vdots \\ (S_{19}, a_{13}, r_{19}, S_{20}) \end{cases}$

run several episodes.
put the record of these episodes
into a memory — experience relay memory
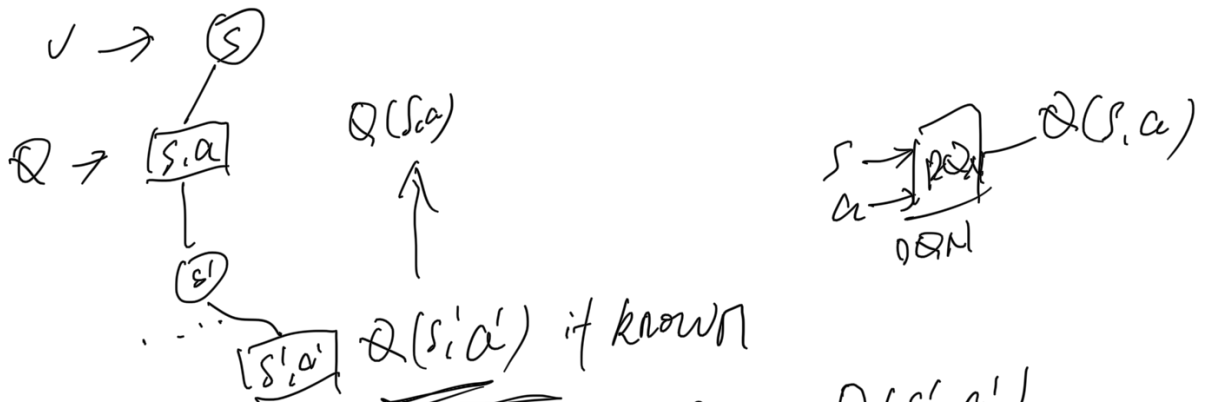remove old episodes, add new episodes.

Memory $\Rightarrow$ training database.

Batch for training, randomly select steps in
in memory, calculate loss

②            $L = ||$ target $-$ current_Network_output $||^2$

$\checkmark \rightarrow$ Ⓢ

$Q \rightarrow$ $\boxed{S, a}$            $Q(S, a)$

$\boxed{S'}$

$\cdots\cdots$ $\boxed{S', a'}$ $Q(S', a')$ if known

$S \rightarrow \boxed{DQN} - Q(S, a)$
$a \rightarrow$
        DQN

target $\underline{Q(S, a)} = \gamma + \lambda \max_{a'} \underline{Q(S', a')}$

        ground truth            current DQN
                $\downarrow$

instead of $Q \sim Q^{*}$,

we do $Q \sim Q\_target \checkmark$

$$L = ||Q(S, a, \omega) - [\gamma + \lambda \max_{a'} Q(S', a', \omega)]||^2$$

            $\uparrow$                        $\uparrow$
            DQN                        DQN

for one state, one action, $\Rightarrow$ one step in the batch
we one $L$, for $n$ steps in the batch.

$l = \sum_i l_i$

③

$$W = W - \gamma \frac{\partial L}{\partial W} \quad , \quad \text{opitimization}$$

$\hookrightarrow$ give a new $Q(s, a)$

---

revisite memory:

At beginning, DQN is useless, because output is random., random episodes are ok,
when you get better DQN, you want to use it to select action.

$\varepsilon$-greedy, with DQN,

stochastic
Prioritization

$\boxed{\begin{array}{c} s, a, r, s' \\ | \\ | \\ | \\ ( \\ ( \end{array}} \quad P(i) = \dfrac{L_i^\alpha}{\sum\limits_{k} L_k^\alpha} \qquad \text{if } \alpha = 0, \text{ uniform}$

---

Double $Q$-learning:

Train two DQNs: $Q_1$, $Q_2$

when decide which action to take in an episode
use $\varepsilon$-greedy with $Q_1 + Q_2$

Calculating Loss.

$$L = (\text{target} - \text{current DQN output})^2$$

$$\text{target} = r + \lambda \max_{a'} (Q(s', a'))$$

with 50/50 chance :

$$\rightarrow L = r + \lambda \max_{a'} Q_1(s', a') - Q_2(s, a)$$

↑ select $a'$, use $Q_2(s', a')$
to select $a'$
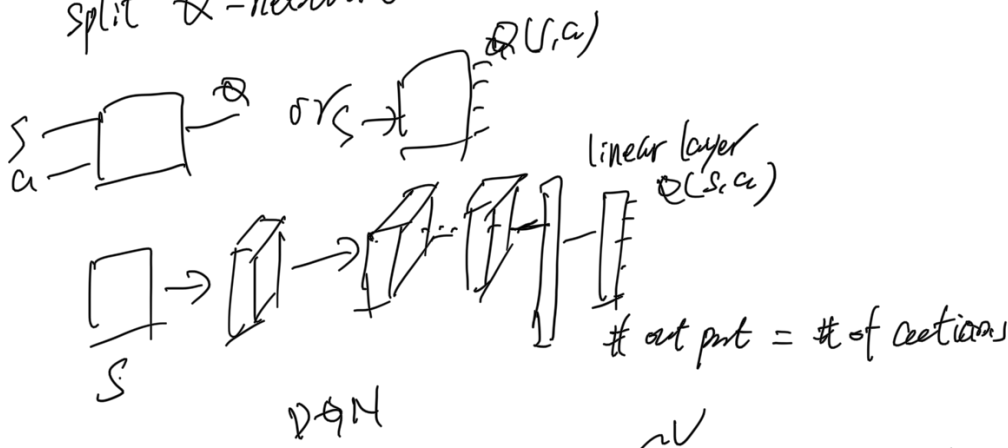
$\rightarrow$ update $Q_2$

another 50/50 chance :

$$\rightarrow L = r + \lambda \max_{a'} Q_2(s', a') - Q_1(s, a)$$

↳ select $a'$, use $Q_1(s', a')$

$\rightarrow$ update $Q_1$

Dueling network.

split $Q$-network into two



DQN

$Q(s, a)$

linear layer
$Q(s, a)$

# of out port = # of actions
$\sim v$

$Q(s, a)$

V — function, Q — function

$$Q(s, a) = V(s) + A(s, a)$$

$$A(s, a) = Q(s, a) - V(s)$$

↑ Advantage function.