# Lecture 13

State - action - reward - state - action (SARSA)

Q - learning:

Go through a number of episodes

    start an episode with an $S$

    Go through the episode step by step

      $Q = Q\text{-new}$ —— give you policy

Sarsa

    based on $Q$, choose an action $a$. —— $\epsilon$-greedy

        base on

        $V = \max_{a} (Q(s, a))$

        or

        random

Take action $a$, get reward, get to a new state $s'$

$$Q_{\_curr} = r + \lambda \underline{V(s')}$$

— Update

$$Q_{\_new}(s, a) = Q(s, a) + \alpha (Q_{\_curr}(s, a) - Q(s, a))$$

if $\alpha = 1$

$s \leftarrow s'$

continue until to the end of the episode.

Sarsa:

Go through a number of episode

start with an S
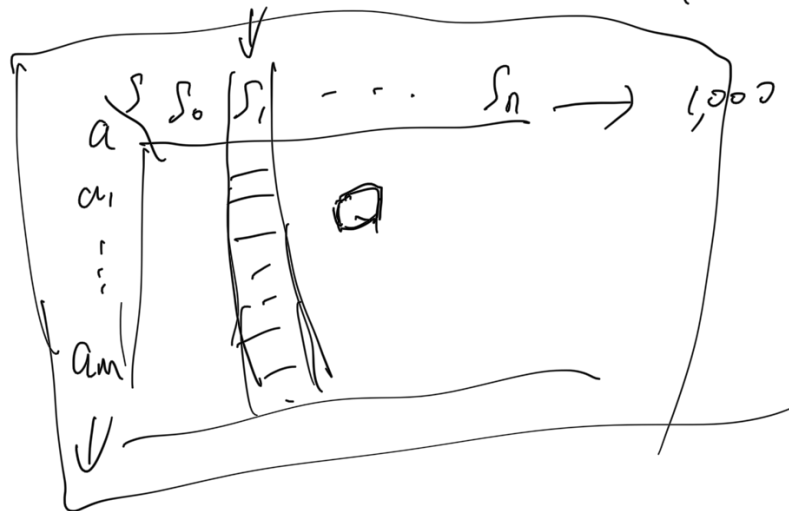
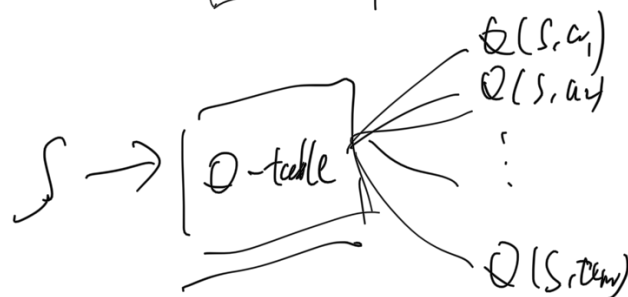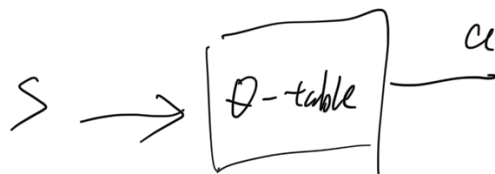→ Choose actions use the current $Q$ ( $\epsilon$-greedy)

then carry out those chosen actions.

take one action at a time, get reward,

reach to a new state $S'$

$$Q_{\text{new}}(S-a) \Leftarrow Q(S,a) + \alpha \left( Q_{\text{current}}(S,a) - Q(S,a) \right)$$

$$\uparrow$$

$$\downarrow$$

$a \left\{ \begin{array}{c} S \; S_0 \; S_1 \quad - \cdots \quad S_n \longrightarrow 1,000 \\ a_1 \\ \vdots \\ a_m \end{array} \right.$

$\boxed{Q}$

$\Downarrow$

Q-learn

$S \longrightarrow \boxed{Q\text{-table}} \quad \rule{1cm}{0pt} a$

$S \longrightarrow \boxed{Q\text{-table}}$
$\quad Q(S,a_1)$
$\quad Q(S,a_2)$
$\quad \vdots$
$\quad Q(S,a_m)$

| X | Y |
|---|---|
| $x_1$ | $y_1$ |
| $x_2$ | $y_2$ |
| . | . |
| . | . |
| $x_n$ | $y_n$ |

$\Rightarrow \quad Y = f(x)$

look up table          regression

$S \longrightarrow$
$C_i \longrightarrow$  DNN $\omega$  $\longrightarrow Q(S, a_i)$

$Q(S, a) \longrightarrow Q(S, a, \omega)$

$\underset{\uparrow}{\phantom{Q}}$ neural network

if you know the ground truth $\hat{Q}(S, a)$

input : $(S, a)_k$  ,  out is one-D $\hat{Q}(S, a)$

| s | a | $Q(S, a)$ |
|---|---|---|
| $X_1$ | $X_2$ | Y |

$L = \left( Q(S, a, \omega) - \hat{Q}(S, a) \right)^2$

$S \longrightarrow \boxed{\omega} \longrightarrow Q$
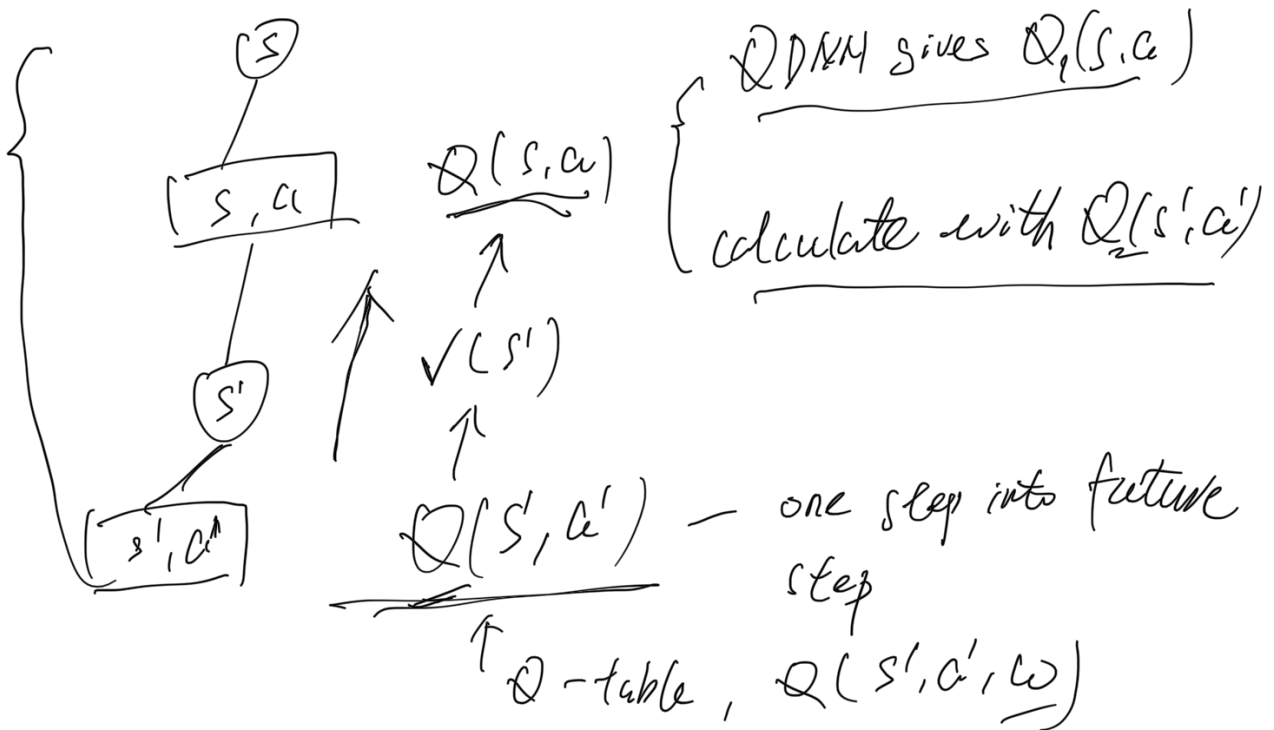$a$

$\dfrac{\partial L}{\partial \omega}$

$Q_0$

Q-train

$$\begin{cases} Q(s,a,\omega) \;\; \text{— current } QDNN \\ Q_{curr}(s,a) = r + \lambda \max_{a'} Q(s',a',\omega) \end{cases}$$

$s'$ is next state after take action $a$.

$$\boxed{s}$$

$$\boxed{s,a}$$

$$\boxed{s'}$$

$$\boxed{s',a'}$$

$\underline{Q(s,a)}$

$\begin{cases} \text{QDNN gives } Q_1(s,a) \\ \text{calculate with } Q_2(s',a') \end{cases}$

$V(s')$

$\underline{Q(s',a')}$ — one step into future step

$\uparrow$ Q-table, $Q(s',a',\omega)$

$s$
$a$ $\boxed{\omega}$ — $Q(s,a)$
$\uparrow$

$s'$
$a'$ $\boxed{\omega}$ — $Q(s',a')$
$\uparrow$

$(s,a,r,s')$

$$L(\omega) = \Big( Q_{curr}(s,a,\omega) - Q(s,a,\omega) \Big)^2$$

$\uparrow$ contain $r$      $\uparrow$ purely from NN

$$\hat{Q}_1 \uparrow \qquad\qquad Q_0$$

$$\frac{\partial L(w)}{\partial w}$$

$$L = \sum_{i=1}^{M} (y_i - \hat{y}_i)^2 \quad , \quad N \text{ samples. for regular DNN}$$

for one $(S_i, a_i, r_i, S_{i+1})$, we have on $L_i$

$$L = \sum_{i=1}^{N} L_i \quad , \quad N \text{ could be very large.}$$

Get all $(S_i, a_i, r_i, S_{i+1})$,

Sample, put into mini-batch,

calculate loss here mini-batch

Double $Q$-learning:

Train two $Q$'s: $Q_1$ and $Q_2$

Do $Q$-learning on both

$$Q_{1\_cur}(S, a) = r + \lambda \max (Q_2(S', a'))$$

alternate

$$Q_{2,\text{corr}}(s,a) = r + \lambda \max\left(Q_1(s',a')\right)$$