# AI

Env      Agent

$S$

$T(S \to S')$     sense $\to$    $\pi$

$R(s)$        goal

     $A$

finite, discrete state - action space

MDP                RL

$\left\{ \begin{array}{l} T(s,a,s') \\ R(s) \end{array} \right\} \to DP$     $\left. \begin{array}{l} T \times \\ R \times \end{array} \right\} \to$

RL $\Big\{$ model-based : $\hat{T}, \hat{R} \} \to DP \left\{ \begin{array}{l} VI \\ PI \end{array} \right.$

model-free :

     Passive RL : evaluate $(\pi) \to (V)$ $\left\{ \begin{array}{l} MC \text{ (Direct sampling)} \\ TD \text{ (temporal-difference)} \\ = Moving\,average \end{array} \right.$

     active RL $(\pi)$

    $\left\{ \begin{array}{l} MC \text{ control} : Q(s,a) \\ TD \text{ control} : Q \end{array} \right.$      $V_{k+1}(s) \leftarrow V_k^t(s) + \alpha [R_k + \gamma V_k(s') - V_k(s)]$

    on-policy :   $\varepsilon$-greedy $\left\{ \begin{array}{l} \varepsilon \,<\, \boxed{\pi} \\ 1-\varepsilon : explore \end{array} \right.$    $\pi$

    off-policy :    $\pi$

          $b$

Q-learning : off-policy, TD-control

SARSA        : on-policy     ,,
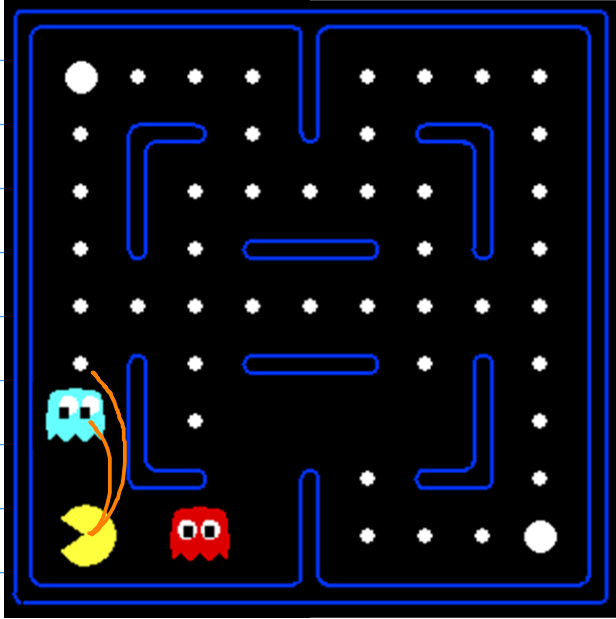
SARSA                 vs         Q-Learning

✓ A from s using $\varepsilon$-greedy Q ;      b: A from s using $\varepsilon$-greedy

. A' from s'    "    "

$Q(s,A) \leftarrow Q(s,A) + \alpha [R + \gamma Q(s',A') - Q(s,A)]$    $\pi$

$Q(s,A) \leftarrow Q(s,A) \; \alpha [R + \gamma \max_a Q(s',a) - Q(s,A)]$

# Approximation Method

$$V = \sum_j W_j \cdot f_j(S)$$

$$Q = \sum_j W_j \cdot f_j(S,a)$$

$$\boxed{Q \leftarrow Q^{old} + \alpha \,(\text{difference})}$$

$$\underline{Q_{targe} - \hat{Q}}$$

$$W_j \leftarrow W_j + \alpha \left(\frac{\partial \ell}{\partial W_j}\right)$$

$$\ell = \frac{1}{2}(Q_t - \hat{Q})^2 \qquad \boxed{(Q_t - \hat{Q}) \cdot \frac{\partial \hat{Q}}{\partial W_j}}$$

$$\boxed{W_j \leftarrow W_j + \alpha \cdot (Q_t - \hat{Q}) \cdot f_j(S,a)}$$

1. <u>non-linear model</u> → ?

2. $\boxed{Q_{true}}$

3. Deep RL

---

## ML

Supervised learning $\ell(y, \hat{y})$

- <u>Parametric model</u> $\quad \hat{y} = f_\theta(x) \quad \{\theta\}$, gradient desc
- non-parametric. $\quad \hat{y} = f(x) \quad \begin{array}{l}\text{dis} \\ \text{ent}\end{array}$

- ) <u>MC</u> → $E(r)$ → $Q(s,a) = G_t$

- ) Bootstrapping $\boxed{r + \gamma \max_{a'} Q_\theta(s', a')}$ → semi-gradient method

$$Q_t$$

$$\left(Q_t - \hat{Q}\right)\frac{\partial \hat{Q}}{\partial \theta}$$

$$W_j \leftarrow W_j + \alpha \frac{\partial \hat{\phi}_W}{\partial W}(s_i, a_i)\left( Q_t - \hat{Q}_W(s_i, a_i)\right)$$

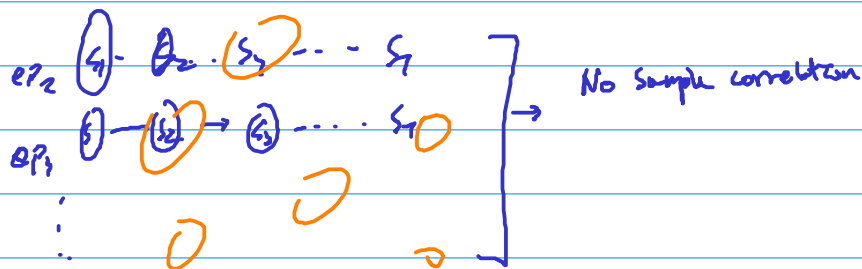$$r(s,a) + \gamma \max_{a'} \hat{\phi}_W(s', a')$$

# Q-learning /w $\boxed{f.A}$

1) sample $(a)$     $s, a \to s', r$    #1 Samples are correlate ✓

2) $Q_t = r(s,a) + \gamma \max_{a'} \hat{Q}_w(s', a')$   #2 off-policy, bootstrapped ✓

3) $w \leftarrow w + \alpha \frac{\partial \hat{Q}(s,w)}{\partial w} \left( Q_t - \hat{Q}(s,a) \right)$   # Semi-gradient descent

## Deep Q-Network 2015

1. F.A $\to$ DNN
2. Off-policy    $\Big\} \to$ never converge!
3. Bootstraping

P1. Correlated samples. $\to$ replay buffer

ep 

$ep_2$ 

$ep_1$    $\to$ No sample correlation

#2,3: Regression target not stable

$w' \leftarrow w$

$\{s_i, a_i, s'_i, r_i\} \to$ RB

$w \leftarrow w + \alpha \sum_i \frac{\partial \hat{Q}_w}{\partial w} \left( r(s_i, a_i) + \gamma \max_{a'} \hat{Q}_{w'}(s'_i, a'_i) - \hat{Q}_w(s_i, a_i) \right)$

       target network    $w' = w$

$\to$ Stable regression

$N = 1000 \sim 10,000$

$w' \leftarrow \tau w' + (1-\tau) w$