# Blackboard RL1 : $\underline{Q - values}$



$Q(s, a_2)$

"green"    $\boxed{a_1}$    $\boxed{a_2}$

$P^{a_2}_{s \to s_4}$

$r_t$

$S_1$   $S_2$   $S_3$   $S_4$

"branching ratio"

Transition probability

$$P^{a_2}_{s \to s_4} = P(s' = s_4 \mid a_2, s)$$

next state

- - - - - - - - - - - - - - - - - - - -

- actual reward at time $t$:    $r_t$

- expected reward for this "branch"

$$R^{a_2}_{s \to s_4} = E(r_t \mid s' = s_4, a_2, s)$$

↑ reward received

↑ end up in $s_4$

↑ take $a_2$

↑ start in $s$

- expected reward for action $a_2$

$$Q(s, a_2) = E(r_t \mid a_2, s)$$

$$= \sum_{s'} P^{a_2}_{s \to s'} \cdot R^{a_2}_{s \to s'}$$

↑ all possible "next states"

# Blackboard RL1-2 = Exercise 1

$Q$ = expected reward $\approx$ empirical mean $r.$ = $\hat{Q}$

$\hat{Q}^{(k-1)}(s,a)$ after $k-1$ trials (playing action $a$)

$$\hat{Q}^{(k-1)}(s,a) = \frac{1}{k-1}\left(r_1 + r_2 + \ldots r_{k-1}\right)$$

↖ 2nd time action $a$

- - - - - - - - - - - - - - - - - - - -

after $k$ trials

$$\hat{Q}^{(k)}(s,a) = \frac{1}{k}\left(r_1 + r_2 + \ldots r_{k-1} + r_k\right)$$

$$= \frac{k-1}{k} \cdot \hat{Q}^{(k-1)}(s,a) + \frac{1}{k}r_k$$

$$= \frac{k}{k}\hat{Q}^{(k-1)}(s,a) + \frac{1}{k}r_k - \frac{1}{k}\hat{Q}^{(k-1)}(s,a)$$

Eq.(1) $\Delta\hat{Q}(s,a) = \hat{Q}^{(k)}(s,a) - \hat{Q}^{(k-1)}(s,a) = \frac{1}{k}\left[r_k - \hat{Q}^{(k-1)}\right]$

$$\Rightarrow \left\|\left| \eta = \frac{1}{k} \right|\right\|$$

theorem (i): if $\quad \underset{\nearrow}{E}\left[\Delta\hat{Q}(s,a)\right] = 0 \qquad\qquad (H)$

then $\quad \underset{\underset{\text{expectation}}{\uparrow}}{E}\left[\hat{Q}(s,a)\right] = \underset{s'}{\sum} P^a_{s\to s'} R^a_{s\to s'}$

proof: $\qquad$ (H) $\qquad$ Eq.(1) of slide
$$\qquad\qquad\qquad \downarrow \qquad\qquad \downarrow$$
$$E\left[\Delta\hat{Q}(s,a)\right] \overset{!}{=} 0 = E\left[r_t - \hat{Q}(s,a)\right]$$

$\underset{\underset{\text{around zero}}{\text{fluctuates}}}{\uparrow} \qquad 0 = E[r_t] - E\left[\hat{Q}(s,a)\right]$

$$0 = \underset{s'}{\sum} P^a_{s\to s'} R^a_{s\to s'} - E\left[\hat{Q}(s,a)\right]$$

(ii) Fluctuations: role of $\eta$ is qualitatively obvious.

$\qquad\qquad\qquad$ smaller $\eta \Rightarrow$ smaller fluctuation

$\uparrow$

$\hat{Q}$ fluctuates around $E\left[Q(s,a)\right] = Q(s,a)$

$\qquad\qquad\qquad\qquad\qquad \underset{\text{expectation}}{\uparrow}$

Blackboard RL4-4 : Exercise 2

update with $\Delta Q(s,a) = 0.2 \cdot [r_t - Q(s,a)]$ (*)

<u>2.1.</u> initialise $Q(s,a_1) = Q(s,a_2) = 0$

$t=1$, action $a_1$; $r_t = 1 \Rightarrow Q(s,a_1) = 0.2$

$t=2$, action $a_2$; $r_t = 0.4 \Rightarrow \boxed{Q(s,a_2) = 0.08}$

<u>2.2.</u> $t=3$, best action $= a_1$; $r_t = 0$

$Q(s,a_1) \leftarrow Q(s,a_1) + 0.2[0 - 0.2]$; $\underline{Q(s,a_1) = 0.16}$

$t=4$, best action $a_1$; $r_t = 0$

$Q(s,a_1) \leftarrow Q(s,a_1) + 0.2[0 - 0.16]$

$\phantom{Q(s,a_1) \leftarrow} 0.16 \quad - \quad 0.032 \qquad \underline{Q(s,a_1) = 0.128}$

$t=5$, best action $a_1$; $r_t = 0$
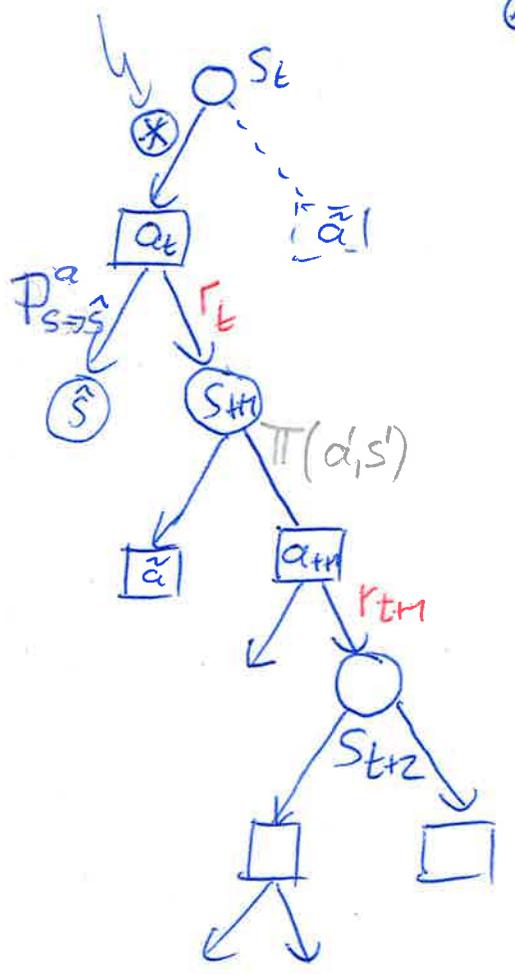
$Q(s,a_1) \leftarrow 0.128 - 0.2 \cdot 0.128$; $\underline{Q(s,a) \cong 0.102}$

$\Rightarrow$ $\underline{a_1}$ remains "best action" for several steps!

<u>2.3</u> actual values

$\left. \begin{array}{l} Q(s,a_1) = 0.25 \\ Q(s,a_2) = 0.30 \end{array} \right\} \Rightarrow \underline{a_2 \text{ is best action}}$

we start here



$P^a_{s \to \hat{s}}$   $r_t$

$\pi(a',s')$

$\circledast$ total reward collected in single trial starting in $s$ with action $a_t$

$$R^{tot}(s_t, a_t) = r_t + \gamma\, r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

$$= r_t + \gamma \left[ r_{t+1} + \gamma\, r_{t+2} + \gamma^2 r_{t+3} + \dots \right]$$

$$= r_t + \gamma \cdot R^{tot}(s_{t+1}, a_{t+1})$$

total reward (single trial) starting from $s' = s_{t+1}$ with $a_{t+1}$

now we look at diagram to calculate expectation

$$E\left( R^{tot}(s_t, a_t) \right) = E\left( r_t + \gamma\, R^{tot}(s_{t+1}, a_{t+1}) \right)$$

$$= \sum_{s'} P^{a_t}_{s \to s'} \left[ R^{a_t}_{s \to s'} + \gamma\, E\left( R^{tot} | s' \right) \right]$$

↑ starting in $s'$

$$= \sum_{s'} P^{a_t}_{s \to s'} \left[ R^{a_t}_{s \to s'} + \gamma \cdot \sum_{a'} \pi(a',s') E R(s',a') \right]$$

$$Q(s_t, a_t) = \sum_{s'} P^{a_t}_{s \to s'} \left[ R^{a_t}_{s \to s'} + \gamma \sum_{a'} \pi(a',s')\, Q^{tot}(s',a') \right]$$

from diagram

$$Q(s,a) \quad \approx \quad r_t \; + \; \gamma \cdot Q(s',a')$$

discount ↓ (above $\gamma$)

$$0 \quad \approx \quad r_t + \gamma \cdot Q(s',a') - Q(s,a)$$

proposed update

$$\Delta Q(s,a) \; = \; \eta \left[ r_t + \gamma \cdot Q(s',a') - Q(s,a) \right]$$

check:

$$\text{if} \quad r_t \; > \; \underbrace{\gamma \cdot Q(s',a') - Q(s,a)}_{} \quad \Rightarrow \quad \text{increase } Q(s,a)$$

↑
actual
reward

expected reward
for this transition

SARSA update

$$\Delta Q(s,a) = \eta \left[ r_t + \gamma\, Q(s',a') - Q(s,a) \right]$$

hypothesis

$$E\left[\Delta Q(s,a)\right] \overset{\downarrow}{=} 0 = \underset{\underset{\text{starting in } s_t \text{ with } a}{\uparrow}}{E}\left[ r_t + \gamma\, Q(s',a') - Q(s,a) \right]$$

$$0 = \sum_{s'} P^a_{s\to s'} \left[ R^a_{s\to s'} + \gamma \sum_{a'} \Pi(s',a')\, Q(s'\,a') \right] - Q(s,a)$$

$\Rightarrow$ Bellman ✓

in order to evaluate expectations:

- look at graph!

- if I am in $s$, all remaining expectations are "given $s$"

- if I am in a branch $(s, a)$ all remaining expectation are given $s$ and $a$