# Blackboard 8.1 : Q - values

"branching ratio"

$s$

$Q(s, a_2)$

"green"   $a_1$   $a_2$

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

$r_t$

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

Transition probability

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

$\uparrow$ next state

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

- expected reward for this "branch"

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

$\uparrow$ reward received

$\uparrow$ end up in $S_4$

$\uparrow$ take $a_2$

$\uparrow$ start in $s$

- expected reward for action $a_2$

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

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

$\uparrow$ all possible "next states"

# Blackboard 8.2 = Exercise 1

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

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

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

$\uparrow$ 2nd time action $a$

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

after $k$ trials

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

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

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

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

$$\Rightarrow \boxed{\eta = \frac{1}{k}}$$

theorem (i):　if　$E\left[\Delta Q(s,a)\right] = 0$　　　　(H)

then　$E\left[Q(s,a)\right] = \sum_{s'} P^a_{s\to s'}\, R^a_{s\to s'}$

↑
expectation

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

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

↑
fluctuates
around zero

$$0 = E\left[r_t\right] - E\left[Q(s,a)\right]$$

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

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

Blackboard 8.4 — Exercise 2    ④

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

2.1. 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 $a_2$; $r_t = 0.4 \Rightarrow \boxed{Q(s,a_2) = 0.08}$

2.2. $t=3$, best action $= a_1$; $r_t = 0$

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

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

$\quad Q(s,a_1) \leftarrow Q(s,a_1) + 0.2[0 - 0.16]$
$\quad\quad\quad 0.16 \quad - 0.032 \quad\quad Q(s,a_1) = 0.128$

$t=5$, best action $a_1$; $r_t = 0$

$\quad Q(s,a_1) \leftarrow 0.128 - 0.2 \cdot 0.128$; $Q(s,a_1) \approx 0.102$

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

2.3 actual values
$\quad\quad Q(s,a_1) = 0.25$
$\quad\quad Q(s,a_2) = 0.30$ $\Big\}$ $\Rightarrow \underline{a_2 \text{ is best action}}$

we start here



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

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

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

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

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

now we look at diagram to calculate expectation

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

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

↑ starting in $s'$

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

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

from diagram                 discount

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

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

proposed update

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

check:

$$\text{if} \quad r_t > \underbrace{\gamma \, Q(s',a') - Q(s,a)}_{\text{expected reward for this transition}} \implies \text{increase } Q(s,a)$$

$\uparrow$
actual
reward

## SARSA update

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

hypothesis

$$E\left[ \Delta Q(s,a) \right] \overset{!}{=} 0 = E \left[ r_t + \gamma Q(s',a') - Q(s,a) \right]$$

starting in $s_t$ with $a$

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

$\Rightarrow$ Bellman $\checkmark$

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 $a$