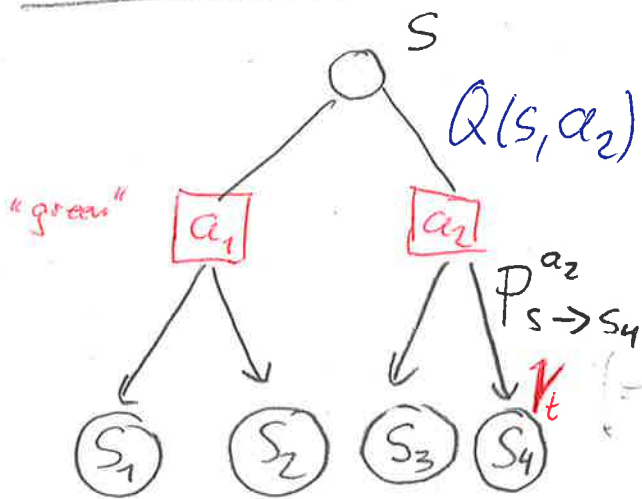


# Blackboard RL1 : Q-values



"branching ratio"

Transition probability

$$P_{s \rightarrow s_4}^{a_2} = P(s' = s_4 | a_2, s)$$

↑  
next state

- actual reward at time t:  $V_t$
- expected reward for this "branch"

$$R_{s \rightarrow s_4}^{a_2} = E(V_t | s' = s_4, a_2, s)$$

↑ reward received      ↑ end up in  $s_4$       ↑ take  $a_2$       ↑ start in  $s$

- expected reward for action  $a_2$

$$Q(s, a_2) = E(V_t | a_2, s)$$

$$= \sum_{s'} P_{s \rightarrow s'}^{a_2} \cdot R_{s \rightarrow s'}^{a_2}$$

↑ all possible "next states"

} (\*Def Q)

- distinguish  $V_t$  from  $R_{s \rightarrow s'}^{a_2}$  ! for example
 

$V_t = 4$	with Prob	$P_r(r=4   s, a, s')$
$V_t = 1$	" "	$P_r(r=1   s, a, s')$
$V_t = 0$	" "	$P_r(r=0   s, a, s')$

Theorem (i): if  $E[\Delta \hat{Q}(s,a)|s,a] = 0$  (H)

then  $E[\hat{Q}(s,a)] = \sum_{s'} P_{s \rightarrow s'}^\alpha R_{s \rightarrow s'}^\alpha$

note: expectation is given (s,a)!

proof: (H) Eq. (1) of slide estimate of Q, given (s,a)

$$E[\Delta \hat{Q}(s,a)|s,a] = 0 = \eta \cdot E[r_t - \hat{Q}(s,a)]$$

updates average to 0 in expectations

$$0 = E[r_t] - E[\hat{Q}(s,a)]$$

$$0 = \sum_{s'} P_{s \rightarrow s'}^\alpha R_{s \rightarrow s'}^\alpha - E[\hat{Q}(s,a)] = Q(s,a)$$

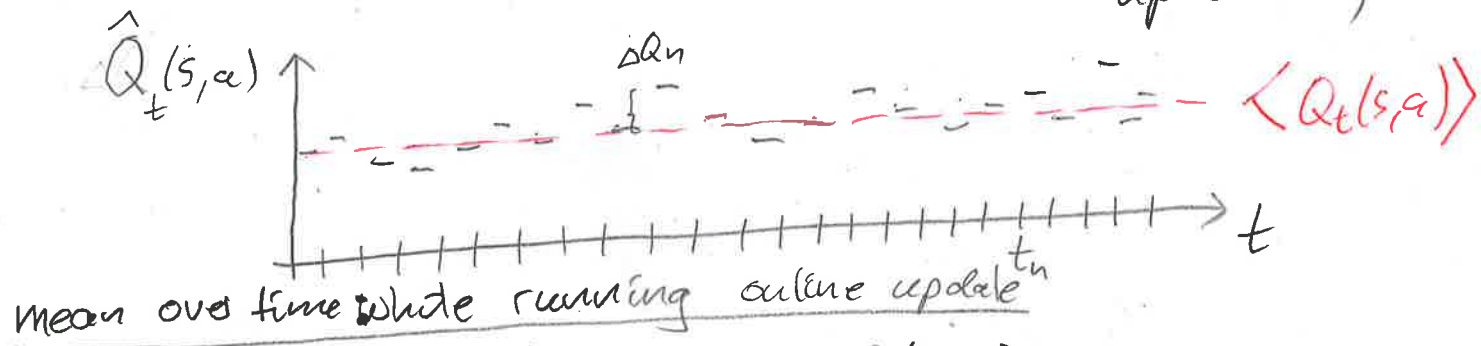
(ii) Fluctuations over time around empirical mean  $\langle \hat{Q} \rangle$ . role of  $\eta$  is qualitatively obvious. smaller  $\eta \Rightarrow$  smaller fluctuation

$\hat{Q}$  fluctuates around  $\langle \hat{Q}(s,a) \rangle = Q(s,a)$   
long-term average expectation

repeat same calculation:

$$\langle \Delta \hat{Q}(s,a) \rangle = 0 = \langle r_t \rangle - \langle \hat{Q}(s,a) \rangle$$

(ii) convergence of mean (perspective of "online" update) ③



$$\langle \Delta \hat{Q}(s, a) | s, a \rangle_{t|s, a} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \Delta Q_n(s, a)$$

update in time step  $n$

theorem (ii)

if  $\langle \Delta \hat{Q}(s, a) | s, a \rangle_{t|s, a} = 0$

then  $\langle \hat{Q}(s, a) \rangle_{t|s, a} = Q(s, a)$

proof:  $\langle \Delta \hat{Q}(s, a) | s, a \rangle_{t|s, a} = \eta \langle (r_t - \hat{Q}_t(s, a)) | s, a \rangle_{t|s, a}$

$$\langle \Delta \hat{Q}(s, a) | s, a \rangle_{t|s, a} = \eta \langle r_t | s, a \rangle_{t|s, a} - \eta \langle Q_t(s, a) \rangle_{t|s, a}$$

↑↑  
automatic conditioning

↓ H

$$0 = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a - \langle Q_t(s, a) \rangle$$

$$\Rightarrow \langle Q_t(s, a) \rangle_t = \sum_s P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a \stackrel{(*\text{def } Q)}{=} Q(s, a)$$

⇒ mean over time  $\langle Q_t \rangle_t = Q$  (exact Q-value)

$Q'$   $Q_t$

see

Blackboard RL1-3 : Stranger Proof (4)

Convergence in expectation

(perspective of "batch" update)

theorem (i) :

if  $E[\Delta \hat{Q}(s,a) | s,a] = 0$  (H)

then  $E[\hat{Q}(s,a) | s,a] = \hat{Q}(s,a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a = Q(s,a)$

proof :

update formul

$\wedge$  = momentary estimate

$E[\Delta \hat{Q}(s,a) | s,a]$

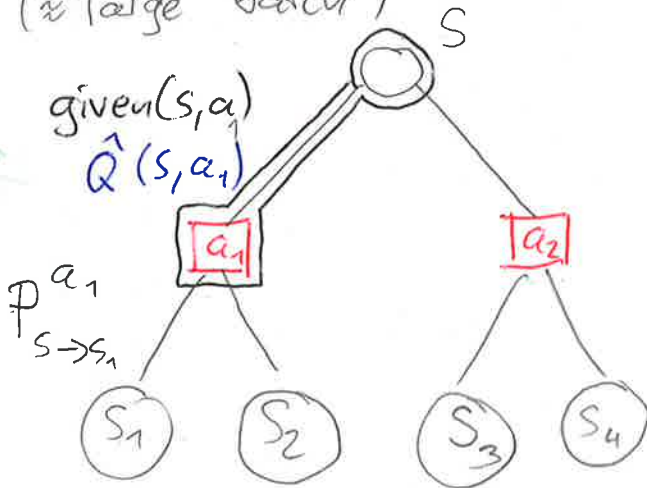
$= E[\underbrace{\eta \cdot (\tau_t - \hat{Q}(s,a))}_{\text{linear}} | s,a]$

possible futures with correct statistical weight (x "large batch")

given s and a

$= \eta \cdot E[\tau_t | s,a] - \eta E[\hat{Q}(s,a) | s,a]$

$= \eta \cdot \sum_{s'} P_{s \rightarrow s'}^a \cdot R_{s \rightarrow s'}^a - \eta E[\hat{Q}(s,a) | s,a]$



$\hat{Q}(s,a)$  is fixed when conditioned on  $(s,a)$   
 $\rightarrow$  drop Expectation sign

with hypothesis H :

$0 = \cancel{\eta} \cdot \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a - \cancel{\eta} E[\hat{Q}(s,a) | s,a]$

$E[\hat{Q}(s,a) | s,a] = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a = Q(s,a)$   
 (\* Def Q)

$\Rightarrow$  estimate  $\hat{Q}$  equal "exact Q"