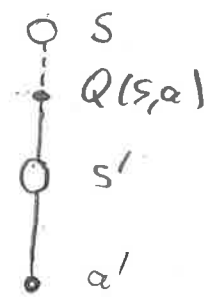
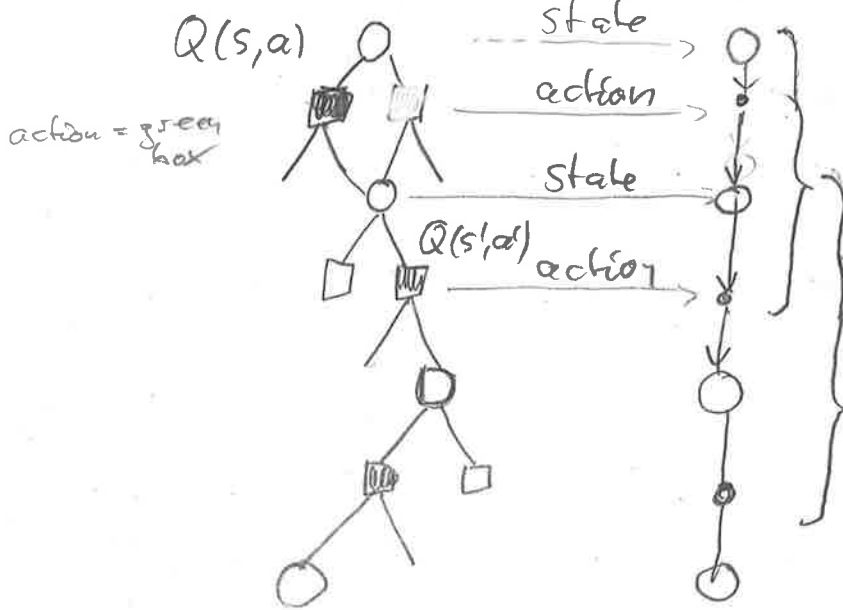


Environment → path → update



$$Q(s,a) \leftarrow Q(s,a) + d \left[\underset{\substack{\uparrow \\ \text{earlier} \\ \text{action}}}{R_t + \gamma Q(s',a')} - Q(s,a) \right]$$

\uparrow
next
action

Consistency of fluctuating online SARSA

right Blackboard

with Bellman equation [from RL 1]

hypothesis: $0 = \langle \Delta \hat{Q}(s,a) \rangle = \underbrace{\eta}_{\text{cut}} \langle r_t + \gamma \hat{Q}(s',a') - \hat{Q}(s,a) \rangle$

$\underbrace{\hspace{10em}}_{\text{shift}}$

$$\langle \hat{Q}(s,a) \rangle \underset{\substack{\uparrow \\ \text{fluctuates}}}{\leftarrow} = \langle r_t + \gamma \cdot \hat{Q}(s',a') \rangle \underset{\substack{\uparrow \\ \text{fluctuates}}}{\leftarrow} \text{temporal average over many "trials" } (N \rightarrow \infty)$$

$$\langle \hat{Q}(s,a) \rangle = \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \langle \pi(s',a') \cdot \hat{Q}(s',a') \rangle]$$

Problem: π^a depends on $Q = \pi(s'|a) \hat{Q}(s',a)$
 ↳ slide (2-2)

Solution: ① if η is small, the fluctuations of \hat{Q} are small and fluctuations of policy π^a are "even smaller"

② consider π^a fixed for small enough Q
 \Rightarrow move π out: $\langle \pi^a \hat{Q} \rangle \sim \pi^a \langle Q \rangle$

$$\langle \hat{Q}(s,a) \rangle = \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi^a(s',a') \langle \hat{Q}(s',a') \rangle]$$

$\langle \hat{Q}(s,a) \rangle = Q(s,a)$ solves Bellman equation $\hat{Q}(s',a') = Q(s',a')$

remarks

① *example of π^a "even smaller": ϵ -greedy
 only rank-order of Q -values matters: best/2nd best/...
 if fluctuations $|\Delta \hat{Q}(s',a')| \ll Q(s',\text{best}) - Q(s',2^{\text{nd best}})$
 then π^a remains stable!

② evaluation of averages - look at graph
 - if in state s' all remaining averages are "given s' "

$$\begin{aligned}
 (1) \quad \hat{Q}_t(s', a') &= \langle \hat{Q}_t(s', a') \rangle_t + \delta Q_t(s', a') \\
 &= \underset{\downarrow \text{defn}}{\bar{q}} + \underset{\downarrow \text{defn}}{\delta q}
 \end{aligned}$$

simplified notation

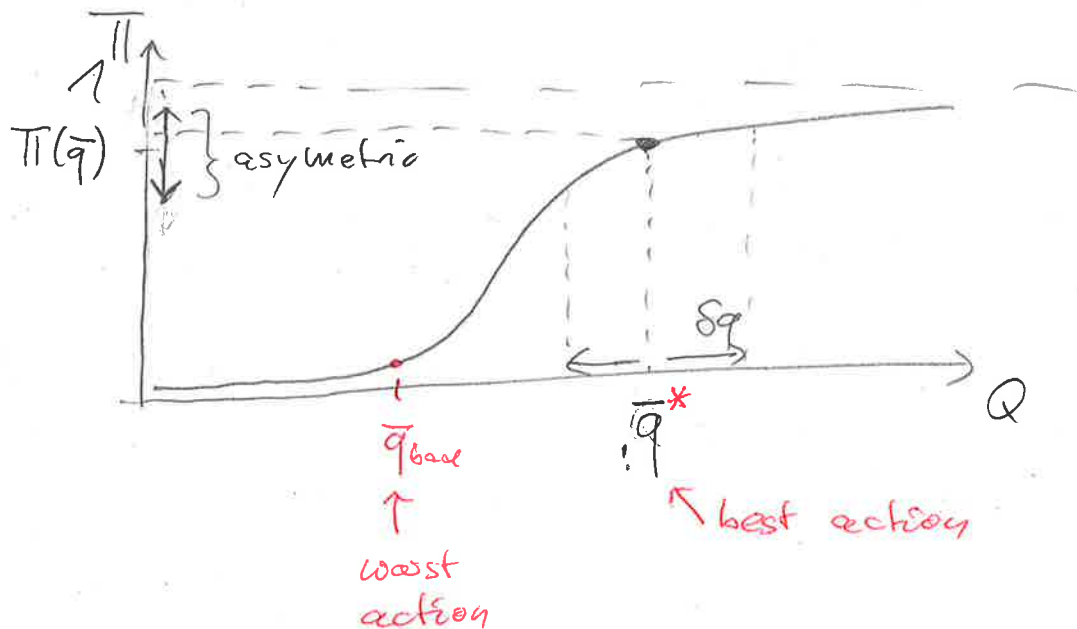
$$\bar{q} := \langle \hat{Q}_t(s', a') \rangle_t$$

$$\delta q := \delta Q_t(s', a') \quad \text{with} \quad \langle \delta q \rangle = 0$$

simplified notation

$$\Pi'(s', a' | \hat{Q}_t(s', a')) = \Pi(\hat{Q}_t) \quad \text{e.g.} \quad \text{softmax}$$

$$(2) \quad \Pi(s', a' | \hat{Q}_t(s', a')) = \Pi(\bar{q}) + \frac{\partial \Pi}{\partial \bar{q}} \cdot \delta q + \frac{1}{2} \frac{\partial^2 \Pi}{\partial \bar{q}^2} \cdot (\delta q)^2$$



We need to evaluate product $\Pi \cdot Q$

$$\Pi(s, a) \hat{Q}(s, a) \cdot \hat{Q}(s, a)$$

$$(1) + (2) = \left[\Pi(\bar{q}) + \delta q \cdot \frac{\partial \Pi}{\partial Q} + \frac{1}{2} (\delta q)^2 \frac{\partial^2 \Pi}{\partial Q^2} \right] \cdot \left[\bar{q} + \delta q \right]$$

collect terms according to $(\delta q)^n$ power

$$= \langle \Pi(\bar{q}) \cdot \bar{q} \rangle$$

drop $\langle \rangle$

$$+ \langle \delta q \rangle \left[\bar{q} \cdot \frac{\partial \Pi}{\partial Q} + \Pi(\bar{q}) \right]$$

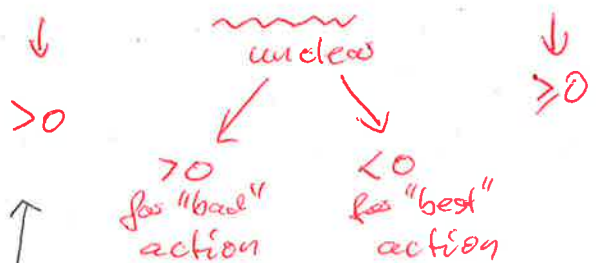
vanishes since $\langle \delta q \rangle = 0$

$$+ \langle (\delta q)^2 \rangle \left[\bar{q} \cdot \frac{1}{2} \frac{\partial^2 \Pi}{\partial Q^2} + \frac{\partial \Pi}{\partial Q} \right]$$

"bias term"

now we average $\langle \rangle$:

$$\langle \Pi(Q_t) \cdot \hat{Q}_t \rangle = \Pi(\bar{q}) \cdot \bar{q} + \langle (\delta q)^2 \rangle \left[\bar{q} \cdot \frac{1}{2} \frac{\partial^2 \Pi}{\partial Q^2} + \frac{\partial \Pi}{\partial Q} \right]$$



term vanishes for $\eta \rightarrow 0$

for finite η : weighted Q-value of "bad" actions is overrated (\rightarrow bias)

but for $\eta \rightarrow 0$ we have $\delta q \rightarrow 0$