

$$Q(s, a) \leftarrow Q(s, a) + d [r_t + \gamma Q(s', a') - Q(s, a)]$$

↑
↑
 earlier next

Exercise 1 (week 9)	trials	Monte Carlo $\langle R(\tilde{s}, \tilde{a}) \rangle$ "average return"	Bootstrap $Q(\tilde{s}, \tilde{a})$ Q from Bellman
(s', a_3) :	2, 4, 8	$\frac{1}{3}(1+1+1)$	$\frac{1}{3}(1+1+1)=1$
(s', a_4) :	1, 3, 6, 7, 9	$\frac{1}{5}(0+0+0+0.5+0.5)$	$\frac{1}{5} \cdot 1$
(s, a_1) :	5, 10	$\frac{1}{2} \cdot (0+0)$	0
(s, a_2) :	1, 9	$\frac{1}{2}(0.2+0.7) = \underline{\underline{0.45}}$; \uparrow	

$$Q(s, a_2) = \frac{1}{2} + \max_a Q(s', a)$$

$$= \underline{\underline{0.2 + 1}}$$

big difference for $Q(s, a_2)$ vs $R(s, a_2)$

$Q(s, a_2)$ much better!

Exercise 1 (continued) : relation online $Q \rightarrow$ batch Q

batch- $Q = Q$ from Bellman

online Q : $Q(s, a) \leftarrow Q(s, a) + \alpha \left[r_t + \underbrace{\max_{a'} \{Q(s', a')\}}_{\substack{\text{fixed,} \\ \text{compressed} \\ \text{knowledge} \\ \text{from previous trials}}} - Q(s, a) \right]$

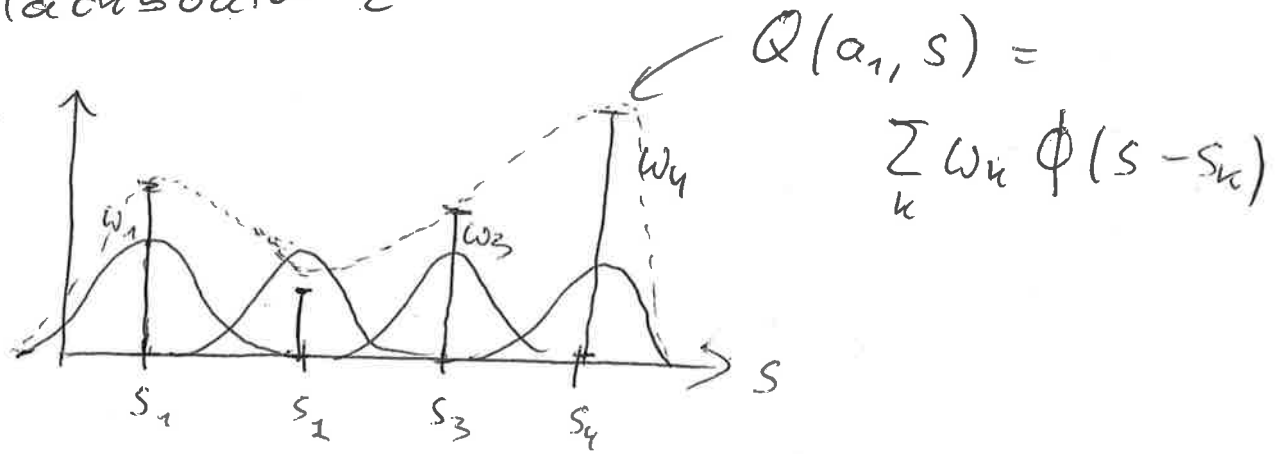
batch Q : n trials starting at s, a , trials $k = 1, \dots, n$

$Q(s, a) \leftarrow Q(s, a) + \frac{1}{n} \left[\sum_{k=1}^n r_t(k) + \max_{a'} \{Q(s', a')\} - Q(s, a) \right]$
initialize with $Q(s, a) = 0$ \downarrow \downarrow
 $=0$ $=0$

$$Q(s, a) \leftarrow \frac{1}{n} \left[\sum_{k=1}^n r_t(k) \right] + \max_{a'} \{Q(s', a')\}$$

$$Q(s, a) \leftarrow \langle r_t \rangle + \max_{a'} \{Q(s', a')\}$$

Blackboard 2



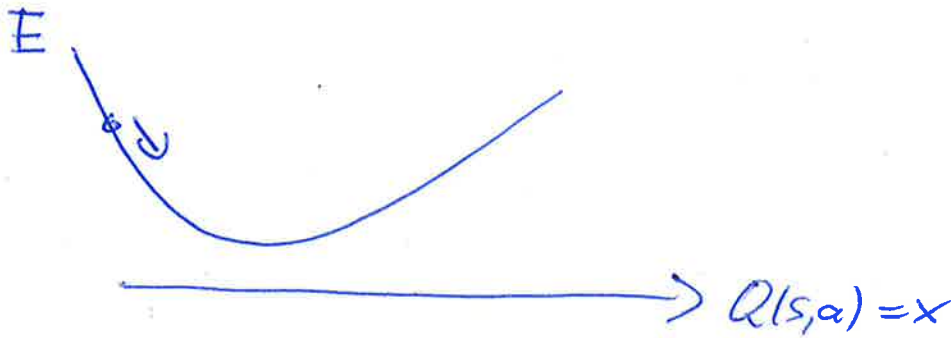
amplitudes

w_1 w_2 w_3 w_4

⇒ smooth function with few parameters

②

$$E = \frac{1}{2} \left[\underbrace{r + Q(s', a')}_{\text{target}} - \underbrace{Q(s, a)}_{\substack{\uparrow \\ \text{variable } x}} \right]^2$$



gradient descent

$$\Delta Q(s, a) = \eta \frac{\partial E}{\partial x} = \underbrace{[r + Q(s', a') - Q(s, a)]}_\delta \cdot \underbrace{\delta_{s, s'} \delta_{a, a'}}_{\text{factor 1}}$$

if action a was chosen in state s

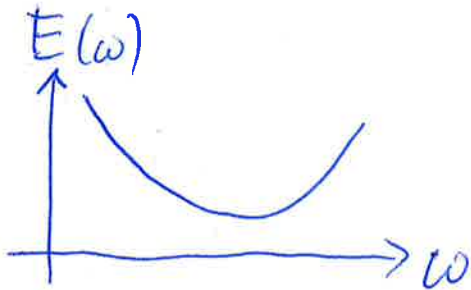
→ SARSA update



Now with

$$Q(s, a) = Q_w(s, a)$$

↑ depends on weights w



$$\Delta w = -\eta \cdot \frac{\partial E}{\partial w} = \underbrace{[r + Q(s', a') - Q_w(s, a)]}_{\text{target}} \cdot \frac{\partial Q_w(s, a)}{\partial w}$$

if: radial basis function + linear output:

$$Q(s, a) = \sum_i w_{ai} \cdot \phi(s - s_i)$$

$$\frac{\partial Q}{\partial w_{bj}} = \phi(s - s_j) \cdot \delta_{ab}$$

↑
only in the action is $a=b$