environment ——→ path ——→ update

$Q(s,a)$

state ———→
action ———→
state ———→
$Q(s',a')$ action ———→

$s$
$Q(s,a)$
$s'$
$a'$

need to know
next action
before you
update $Q(s,a)$
earlier

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

earlier
action

next
action

# RL2, Exercise 1a + 1b          Blackboard 2

| $(\hat{\hat{s}},\hat{\hat{a}})$ pair | encountered in trial | Monte Carlo average return $\langle R(\hat{s},\hat{a})\rangle$ | Bootstrap Batch-Q from Bellman |
|---|---|---|---|
| $(s', a_3^{\vdots})$ | 2, 4, 8 | $\frac{1}{3}[1+1+1]=1$ | $1$ |
| $(s', a_4)$ | 1, 3, 6, 7, 9 | $\frac{1}{5}[0+0+0+0.5+0.5]=\frac{1}{5}$ | $\frac{1}{5}$ |
| $(s, a_1)$ | 5, 10 | $\frac{1}{2}[0+0]=0$ | $0$ |
| $(s, a_2)$ | 1, 9 | $\frac{1}{2}[0.2+0.7]=\underline{0.45}$ | $\langle r_t\rangle + \max_{a'} Q(s',a')$ |

$$\downarrow \qquad\qquad \downarrow$$
$$0.2 + \quad 1 = \underline{\underline{1.2}}$$

---

with same number of trials  Bootstrap/Bellman/Batch-Q yields
much better estimate than  Monte-Carlo!

---

<u>Batch - Q</u> = Q from Bellman (without knowledge of branching ratio)

<u>online Q-learning</u>

$$Q(s,a) \leftarrow Q(s,a) + \eta \cdot \left[ r_t + \max_{a'}\{Q(s',a')\} - \gamma Q(s,a)\right]$$

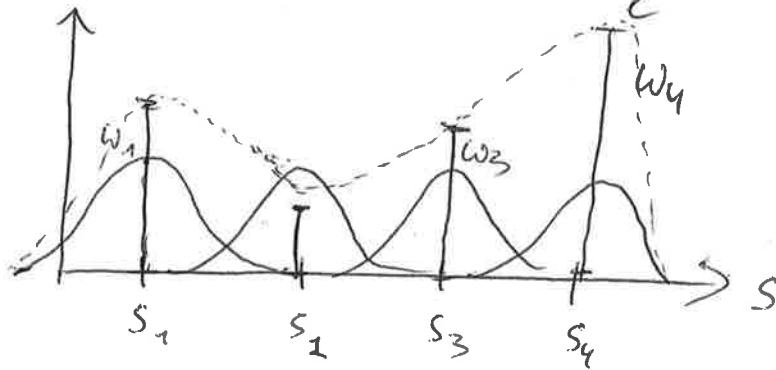<span style="color:red">* compressed knowledge from previous trials</span>

<u>batch - Q</u>:   n trials $(1 \le k \le n)$ starting at $(s,a)$

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

initialize: $Q=0$

$$\downarrow_{0} \qquad\qquad \underset{\text{average}}{\downarrow} \qquad\qquad \underset{\text{* compressed knowledge from states close to target}}{\qquad} \quad \downarrow_0$$

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

$$Q(a_1, s) = \sum_k \omega_k \phi(s - s_k)$$



amplitudes

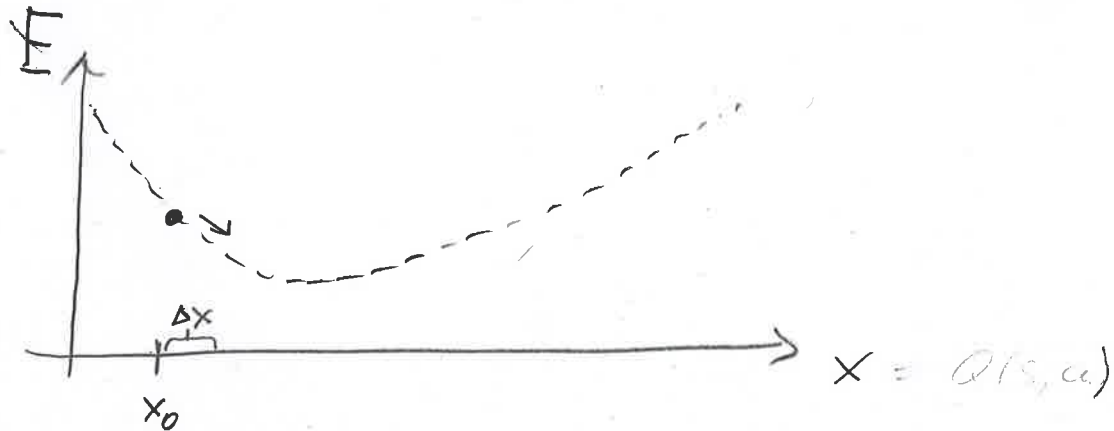$$\omega_1 \quad \omega_2 \quad \omega_3 \quad \omega_4$$

$\Rightarrow$ smooth function with few parameters

## error

$$E = \frac{1}{2}\left[ \underbrace{r + \gamma Q(s',a')}_{\text{target}} - \underbrace{Q(s,a)}_{\substack{\text{depends on}\\ \text{parameter } x \\ (\text{the weights } \omega_1, \omega_2, \ldots)}} \right]^2$$

minimize error by gradient descent

E

$x = Q(s,a)$

$\Delta x$

$x_0$

$$\Delta x_i = -\eta \cdot \frac{\partial E}{\partial x} = +\eta \cdot \left[ r + \gamma Q(s',a') - Q(s,a) \right] \frac{\partial Q(s,a)}{\partial x}$$