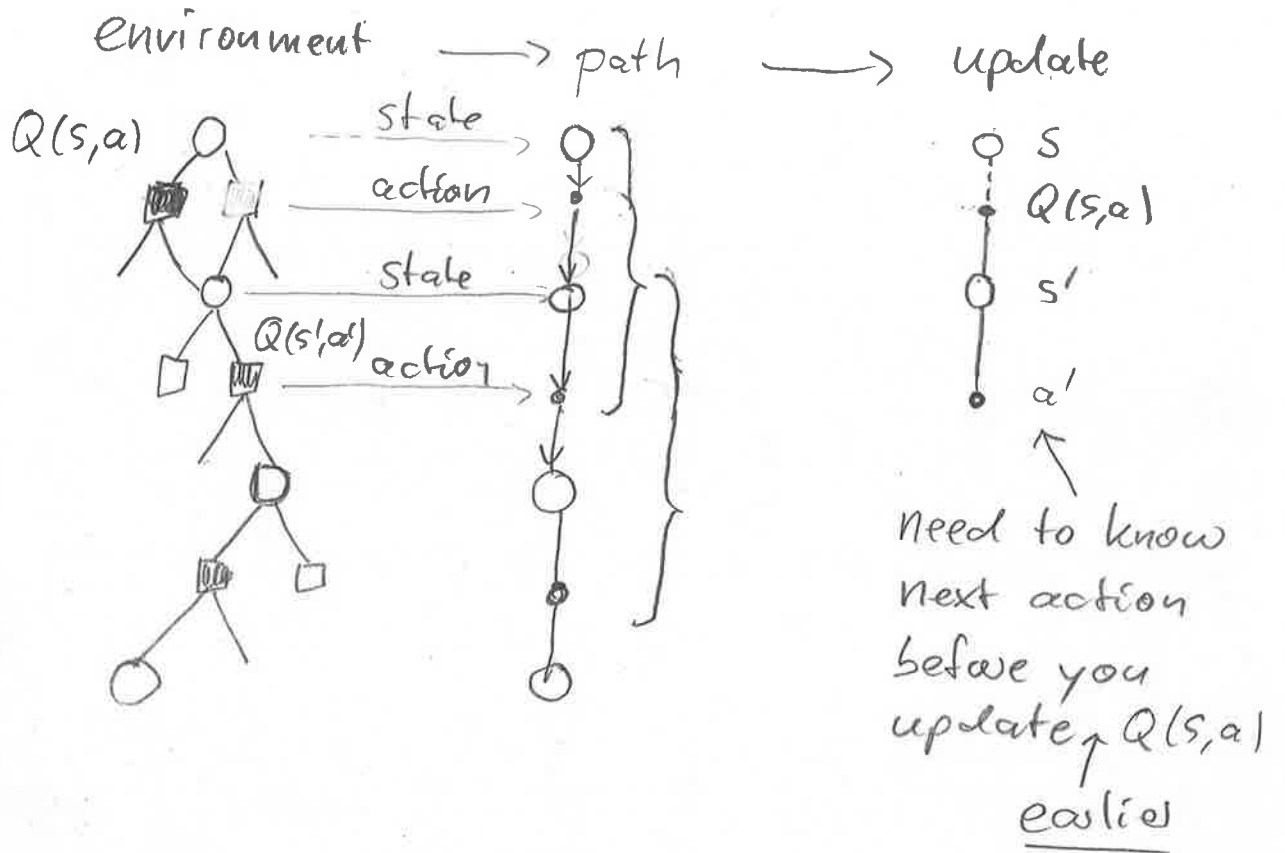


Q-Learning - RL2 ;

Blackboard1: Backup Diagram



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

↑ earlier action ↑ next action

RL2, Exercise 1a + 1b

Blackboard 2

(\hat{s}, \hat{a}) pair	encountered in trial	Monte Carlo average return $\langle R(\hat{s}, \hat{a}) \rangle$	Bootstrap Batch-SARSA from Bellman
(s', a_3)	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	$\frac{1}{2}[0.2+0.7]=0.45$ ↑ only 2 trials are used	$\langle r_t \rangle + \gamma \sum_{a'} \pi(s', a') Q(s', a')$ ↓ $0.2 + \frac{1}{2}(\frac{1}{5} + 1) = 0.8$

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

expected Batch-SARSA = Q from Bellman (without knowledge of branching ratio)

online SARSA - algo

$$Q(s, a) \leftarrow Q(s, a) + \gamma \cdot [r_t + \gamma \cdot \underbrace{\sum_{a'} \pi(s', a') Q(s', a')}_{\text{* compressed knowledge from previous trials}} - Q(s, a)]$$

expected batch-SARSA: n trials ($1 \leq k \leq n$) starting at (s, a)

$$Q(s, a) \leftarrow Q(s, a) + \underbrace{\frac{1}{n} \left[\sum_{k=1}^n r_t(k) \right]}_{\text{average}} + \gamma \underbrace{\sum_{a'} \pi(s', a') Q(s', a')}_{\text{* compressed knowledge from states close to target}} - \gamma Q(s, a)$$

initialize: $Q=0$

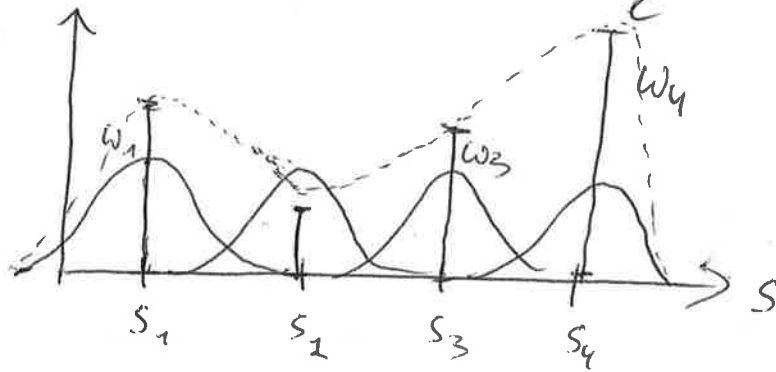
↓
0

since only one s'_j dropped

$$Q(s, a) \leftarrow \langle r_t \rangle + \gamma \left[\sum_{a'} \pi(s', a') Q(s', a') \right]$$

↑ here = 1 0.5

Blackboard 3, RL2



$$Q(a_2, s) = \sum_k w_k \phi(s - s_k)$$

amplitudes

$w_1 \quad w_2 \quad w_3 \quad w_4$

⇒ smooth function with few parameters

$$Q(a_2, s) : \quad \underset{\uparrow}{w_{21}}, \quad \underset{\uparrow}{w_{22}}, \quad \underset{\uparrow}{w_{23}}, \quad \underset{\uparrow}{w_{24}}$$

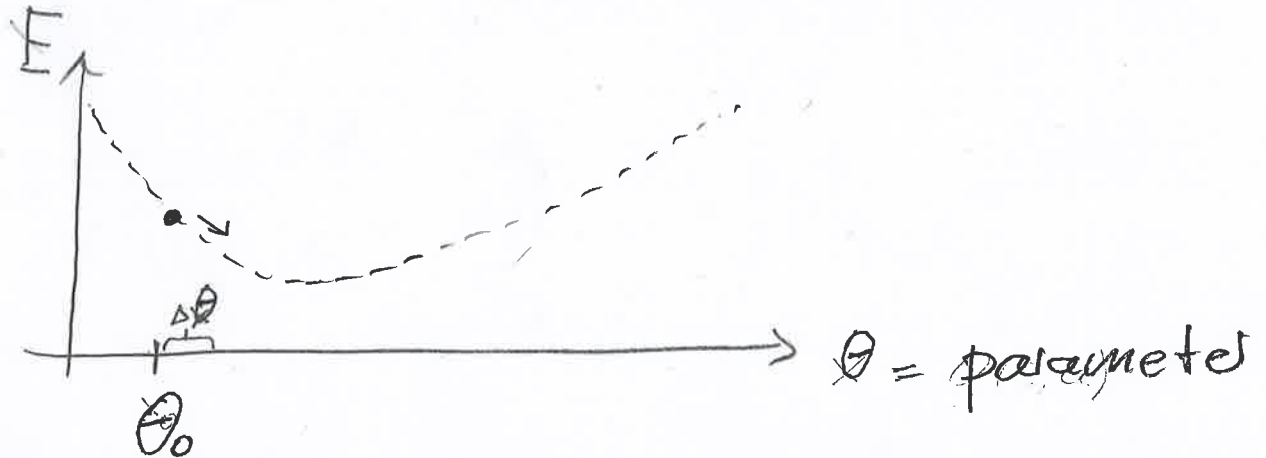
Blackboard 4 - RL 2 :

Loss function

error (loss function)

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

minimize error by gradient descent



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