Unsupervised and Reinforcement Learning
In Neural Networks
Fall 2012
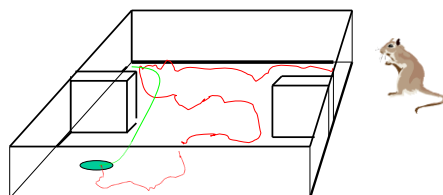
## Reinforcement learning

second lecture on RL:
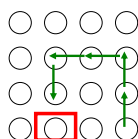- review from previous week
- on-policy (SARSA) vs off-policy (Q-learning)
- eligibility traces

*http://moodle.epfl.ch/*
Wulfram Gerstner
http://lcn1.epfl.ch/

---

**Introduction to reinforcement learning**



LCN

---

## Theory of Reinforcement learning



$s=state$

$a$   $a=action$

$R^a_{s\rightarrow s'} = reward$

$s'=new\ state$

---

**Theory of reinforcement learning – optimal strategy**



$r=R^a_{s\rightarrow s'}$

$Q(s,a1)$   $Q(s,a2)$

$a1$   $a2$

$Q(s,a)= \sum_{s'} P^a_{s\rightarrow s'} R^a_{s\rightarrow s'}$

*Optimal strategy:*
*- take action a\* with*
   $Q(s,a^*)>Q(s,a_j)$
     *other actions*

*But: Q values not known*
  *(probabilities not known)*

  → *Estimate Q values*

Iterative Update
$\Delta Q(s,a)=\eta\ [r-Q(s,a)]$

---

## Reward-based Action Learning

**Exploration versus exploitation**

*Problem: correct Q values not known*
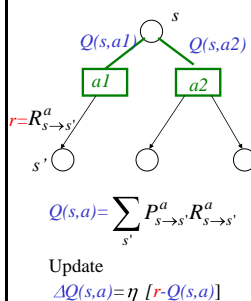 *(since reward probabilities are not known)*

**Exploration versus exploitation**

*Explore so as to estimate reward probabities*

*Take action which looks optimal, so as to maximize reward*

---

## Reward-based Action Learning

**Exploration versus exploitation**



$r=R^a_{s\rightarrow s'}$

$Q(s,a1)$   $Q(s,a2)$

$a1$   $a2$

$Q(s,a)= \sum_{s'} P^a_{s\rightarrow s'} R^a_{s\rightarrow s'}$

Update
$\Delta Q(s,a)=\eta\ [r-Q(s,a)]$

*Problem: correct Q values not known*
*Greedy strategy:*
*- take action a\* which looks best*
   $Q(s,a^*)>Q(s,a_j)]$

$\mathcal{E}$*-greedy strategy:*
*- take action a\* which looks best*
  *with prob*  $P = 1-\varepsilon$

*Softmax strategy:*
$P(action\ a\mid s)=\dfrac{\exp[Q(s,a)]}{\sum_{a'}\exp[Q(s,a')]}$

*Optimistic greedy:*
*start with Q values that are too big*

**Slide 1 (SARSA algo)**

$$Q(s,a) = \sum_{s'} P_{s \to s'}^a \left[ R_{s \to s'}^a + \gamma \sum_{a'} \pi(s',a') Q(s',a') \right]$$

*Bellman equation*

$Q(s,a1)$

$r = R_{s \to s'}^a$

$Q(s',a1)$

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma Q(s',a'))]$

**SARSA algo**

**Slide 2 (SARSA algo)**

**SARSA algo**

*Initialise Q values*
*Start from initial state s*
1) Being in state **s** choose action **a** according to policy
2) Observe reward **r** and next state **s'**
3) Choose action **a'** in state s' according to policy
4) Update

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma Q(s',a'))]$

5) s' → s; a' → a
6) Goto 1)

$Q(s,a1)$  $r = R_{s \to s'}^a$  $Q(s',a1)$

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma Q(s',a'))]$

**Slide 3 (Reward-based Action Learning)**

**Reward-based Action Learning**

$Q(s,a1)$

$r = R_{s \to s'}^a$

$s'$

$$Q(s,a) = \sum_{s'} P_{s \to s'}^a R_{s \to s'}^a$$

Update

$\Delta Q(s,a) = \eta \ [r-Q(s,a)]$

$Q(s,a1)$

$r = R_{s \to s'}^a$

$Q(s',a1)$

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma Q(s',a'))]$

*Small learning rate, algos converge in expectation to the correct solution*

**Slide 4 (Reinforcement learning)**

**Reinforcement learning**

second lecture on RL:

√ - review from previous week

→ - on-policy (SARSA) vs off-policy (Q-learning)

 - eligibility traces

*http://moodle.epfl.ch/*

Wulfram Gerstner

http://lcn1.epfl.ch/

**Slide 5 (Exercise from last week: Bellman equation)**

**Exercise from last week: Bellman equation**

$Q(s,a1)$  $Q(s,a2)$

a1   a2

$r=1$   $r=0$

$s'$

1  2  1  2  0  6

a2    *a=action*

$R_{s \to s'}^a$ = *reward*

 *s'=new state*

*Calculated Q(s,a1) and Q(s,a2)*

a) *All actions equal probability*
b) *Optimal action*

$$Q(s,a) = \sum_{s'} P_{s \to s'}^a \left[ R_{s \to s'}^a + \gamma \sum_{a'} \pi(s',a') Q(s',a') \right]$$

**Slide 6**

$$Q(s,a) = \sum_{s'} P_{s \to s'}^a \left[ R_{s \to s'}^a + \gamma \sum_{a'} \pi(s',a') Q(s',a') \right]$$

*Bellman equations*   *Resulting Q values are different!*

$$Q(s,a) = \sum_{s'} P_{s \to s'}^a \left[ R_{s \to s'}^a + \gamma \max_a Q(s',a') \right]$$

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma Q(s',a'))]$

**SARSA algo**

$\Delta Q(s,a) = \eta \ [r-(Q(s,a)-\gamma \ max_a Q(s',a'))]$

**Q-learning algo**

## Sarsa – on policy temporal difference algo

Initialize $Q(s, a)$ arbitrarily
Repeat (for each episode):
  Initialize $s$
  Choose $a$ from $s$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
  Repeat (for each step of episode):
    Take action $a$, observe $r, s'$
    Choose $a'$ from $s'$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
    $Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma\, Q(s', a') - Q(s, a) \right]$
    $s \leftarrow s'; a \leftarrow a';$
  until $s$ is terminal

*Action for update as used for trajectory*

**Figure 6.9** Sarsa: An on-policy TD control algorithm

$\Delta Q(s,a) = \eta \ [r - (Q(s,a) \gamma\, Q(s',a'))]$

*From: Reinforcement Learning, Sutton and Barto 1998*

---

## Q-learning – off-policy temporal difference algo

Initialize $Q(s, a)$ arbitrarily
Repeat (for each episode):
  Initialize $s$
  Repeat (for each step of episode):
    Choose $a$ from $s$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
    Take action $a$, observe $r, s'$
    $Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$
    $s \leftarrow s';$
  until $s$ is terminal

*Policy used for Action selection*

*Greedy used for update*

**Figure 6.12** Q-learning: An off-policy TD control algorithm.

$$\Delta Q(s,a) = \alpha \cdot \left[ r(a) + \gamma \cdot \max_{s'} Q(s', a) - Q(s,a) \right]$$

*From: Reinforcement Learning, Sutton and Barto 1998*

---

# Reinforcement learning

second lecture on RL:
✓ - review from previous week
✓ - on-policy (SARSA) vs off-policy (Q-learning)
→ - eligibility traces

*http://moodle.epfl.ch/*
Wulfram Gerstner
http://lcn1.epfl.ch/

---

## Exercise from last week

• *Update of Q values in SARSA*

$\Delta Q(s,a) = \eta \ [r - (Q(s,a) - Q(s',a'))]$

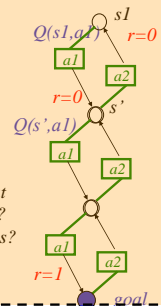• *policy for action choice:*
  *Pick most often action*
$a_t^* = \arg\max_a Q_a(s, a)$

  *Consider a linear sequence*
  *of states. Reward only at goal.*
  *Actions are up or down.*
  *a) Initialise Q values at 0. Start at top. How do Q values develop?*
  *b) Q values after 3 complete trials?*

*goal*

$Q(s1,a1)$   s1   r=0
a1   a2
r=0   s'
$Q(s',a1)$
a1   a2
a1   a2
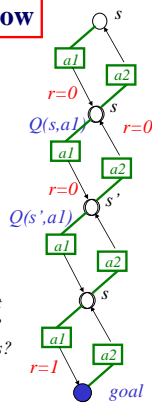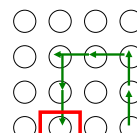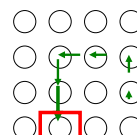r=1   goal



---

## Problem: learning is slow

**- Slow diffusion of information across several states**

$a_t^* = \arg\max_a Q_a(s, a)$

*Consider a linear sequence*
*of states. Reward only at goal.*
*Actions are up or down.*
*a) Initialise Q values at 0. Start at top. How do Q values develop?*
*b) Q values after 3 complete trials?*

s
a1   a2
r=0   s
$Q(s,a1)$   r=0
a1   a2
r=0   s'
$Q(s',a1)$
a1   a2
s
a1   a2
r=1   goal

*goal*



---

## Eligibitiy trace

    ○ s=state
    ⬅a   a=action
    $R_{s \to s'}^a = reward$
    ○ s'=new state

*Sarsa*    $\Delta Q(s,a) = \eta \ [r - (Q(s,a) \gamma\, Q(s',a'))]$

*Idea: at the moment of reward update also previous action values along trajectories*



3

## Slide 1: Eligibitiy trace - algo

**Eligibitiy trace - algo**

*Sarsa(0)*    $\Delta Q(s,a) = \eta \ [r - (Q(s,a)\gamma Q(s',a'))]$

$$\delta_t = r_{t+1} - Q_a(s,a) - \gamma \cdot Q_a(s',a')$$

*Sarsa(λ)*    $\Delta Q_{aj} = \eta \ \delta_t \ e_{aj}$

$if \ a = action \ taken \ in \ state \ j$

$e_{aj}(t) = \gamma\lambda \, e_{aj}(t - \Delta t) + 1$

$else$

$e_{aj}(t) = \gamma\lambda \, e_{aj}(t - \Delta t)$

## Slide 2: Elibitility trace – Sarsa()

**Elibitility trace – Sarsa()**

Initialize $Q(s, a)$ arbitrarily
Repeat (for each episode):
    Initialize $s, a$ and $e(s,a)=0$
    Repeat (for each step of episode):
        Take action $a$, observe $r, s'$
        Choose $a'$ from $s'$ using policy derived from $Q$ (e.g., $\epsilon$-greedy
        $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$
        $e(s, a) \leftarrow e(s, a) + 1$
        For all $s, a$:
            $Q(s, a) \leftarrow Q(s, a) + \alpha\delta e(s, a)$
            $e(s, a) \leftarrow \gamma\lambda \, e(s, a)$
        $s \leftarrow s'; a \leftarrow a'$
    until $s$ is terminal

**Figure 7.11** Tabular Sarsa(λ).

*From: Reinforcement Learnin
Sutton and Barto 1998*

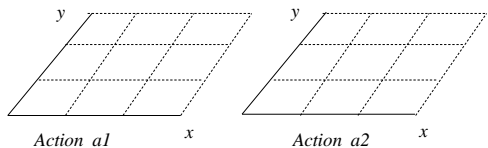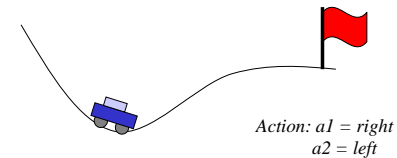## Slide 3: Reinforcement learning

# Reinforcement learning

second lecture on RL:
√ - review from previous week
√ - on-policy (SARSA) vs off-policy (Q-learning)
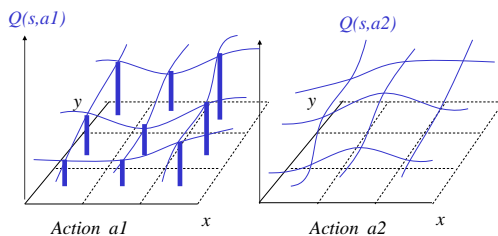√ - eligibility traces
→ - continuous state space

*http://moodle.epfl.ch/*
Wulfram Gerstner
http://lcn1.epfl.ch/

## Slide 4: Continuous states – example

**Continuous states – example**



*Action: a1 = right
a2 = left*

*y* ... *x*    *Action a1*
*y* ... *x*    *Action a2*

## Slide 5: Continuous states – example

**Continuous states – example**

*Blackboard:*

$Q(s,a1)$

$Q(s,a2)$

*Action a1*    *x*
*Action a2*    *x*

## Slide 6: Continuous states

**Continuous states**

$$Q(s,a) = \sum_j w_{aj} r_j(s)$$

$$Q(s,a) = \sum_j w_{aj} \Phi(s - s_j)$$

$Q(s, n)$    $Q(s, s)$

$w_{aj} = w_{n1}$

$r_1(s)$ ... $r_s(s)$

**Spatial Representation**

*Exercise:*  $\dfrac{dQ(s',a')}{dw_{aj}}$

4

## Slide 1

$$Q(s,a) = \sum_{s'} P^a_{s \to s'} \left[ R^a_{s \to s'} + \gamma \sum_{a'} \pi(s',a') Q(s',a') \right]$$

*Bellman equation*

*Q(s,a1)*

a1    a2

$r = R^a_{s \to s'}$

*Q(s',a1)*

a1    a2

$Q(s,a) = r + \gamma\, Q(s',a')$

## Slide 2 — Eligibility and continous: TD( )

- **TD error in SARSA**

$$\delta_t = R_{t+1} - Q_a(s,a) - \gamma \cdot Q_{a'}(s',a')$$

- **policy for action choice:**
  Pick most often action

$$a_t^* = \arg\max_a Q_a(s,a)$$

- **Eligibility trace**

$$e_{aj}(t) = \gamma\lambda\, e_{aj}(t - \Delta t) + \begin{cases} r_j & \text{if } a = action\ taken \\ 0 \end{cases}$$

*memory*

- **Function approximation**

$$Q_w(s,a) = \sum_{i=1}^{n} w_i^a \cdot r_i(s)$$

- **Weight update**

$$\Delta w_{aj} = \eta\, \delta_t\, e_{aj}$$

## Slide 3 — Eligibility and continous: TD( )

- **TD error in SARSA**

$$\delta_t = R_{t+1} - Q_a(s,a) - \gamma \cdot Q_{a'}(s',a')$$

- **Function approximation**

$$Q_w(s,a) = \sum_{i=1}^{n} w_i^a \cdot r_i(s)$$

- **Synaptic update**

$$\Delta w_{aj} = \eta\, \delta_t\, e_{aj}$$

- **Eligibility trace** (*memory at synapse*):

$$e_{aj}(t) = \gamma\lambda\, e_{aj}(t - \Delta t) + \begin{cases} r_j \bullet 1 & \text{if } a = a'\ (= action\ taken) \\ r_j \bullet 0 & else \end{cases}$$

$$e_{aj}(t) = \gamma\lambda\, e_{aj}(t - \Delta t) + r_j\, \delta_{a,a'}$$

## Slide 4 — Reward-based Action Learning

Connection reinforced if
action a at state s  successful

Success signal

**Learning Rule**

$$\Delta w_{aj} = \eta\, \delta_t\, e_{aj}$$

Q(s, n)    Q(s, s)

$w^n$

$r_1(s)$    $r_n(s)$

**Spatial Representation**

## Slide 5 — Example – mountain car



214    Generalization and Function Approximation

MOUNTAIN CAR    Goal

reward
r = -1
except at goal

Step 428    Episode 12

Episode 104    Episode 1000    Episode 9000

Figure 8.10  The mountain-car task (upper left panel) and the cost-to-go function ($\approx \max_a Q_t(s,a)$) learned during one run.
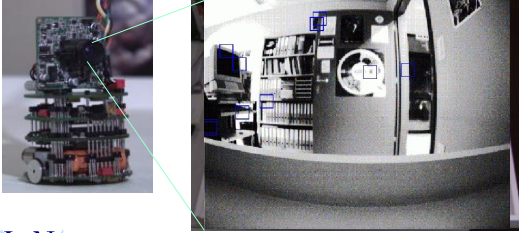
Example 8.2: Mountain-Car Task  Consider the task of driving an underpowered

## Slide 6

**Biological Principles of Learning: spatial learning**



LCN

5

## Validating the Model

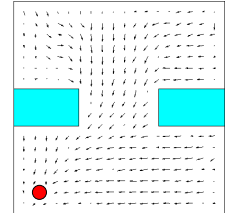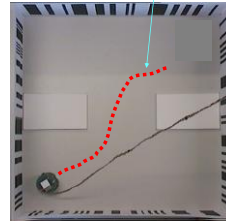**The KHEPERA mobile miniature robot**

**Experimental arena**



**422 x 316 pixels**

LCN

## Open-field Navigation Experiments

*robot trajectory*

*(Biol. Cybern., 2000)*



**Navigation map
after 20 training trials**

LCN

---



11.3  The Acrobot

goal: Raise tip above line

torque applied here

tip

Figure 11.4  The acrobot.

Steps per episode (log scale)

median of 10 runs

typical single run

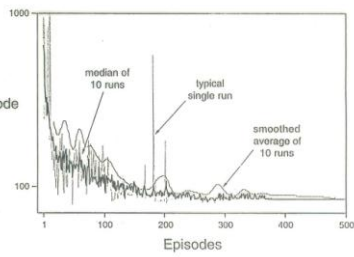smoothed average of 10 runs

Episodes

**Figure 11.6**  Learning curves for Sarsa($\lambda$) on the acrobot task.
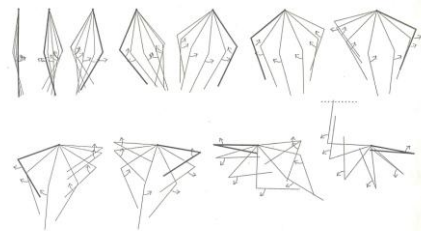
---



274     Case Studies

**Figure 11.7**  A typical learned behavior of the acrobot. Each group is a series of consecutive positions, the thicker line being the first. The arrow indicates the torque applied at the second joint.

---

*The end*