

# Sistemi Intelligenti Reinforcement Learning: Q-learning

Alberto Borghese

Università degli Studi di Milano  
Laboratorio di Sistemi Intelligenti Applicati (AIS-Lab)  
Dipartimento di Scienze dell'Informazione  
[borghese@dsi.unimi.it](mailto:borghese@dsi.unimi.it)



A.A. 2006-2007

1/34

<http://homes.dsi.unimi.it/~borghese/>



## 1° annual RL competition @ NIPS



<http://rlai.cs.ualberta.ca/RLAI/rlc.html>

A.A. 2006-2007

2/34

<http://homes.dsi.unimi.it/~borghese/>



# Sommario



Value e Q functions

SARSA

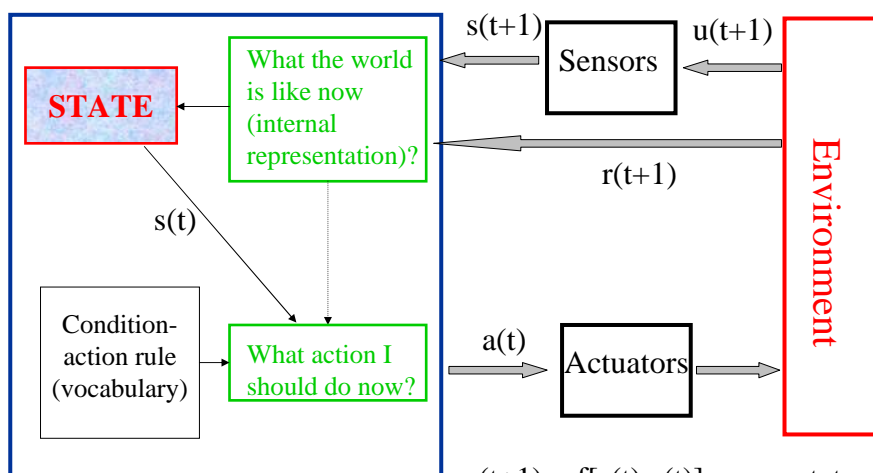
Q-learning



# Schematic diagram of an agent



Agent



s, stato; a, uscita

$$s(t+1) = f[s(t), a(t)]$$

$\pi(s, a)$  è la policy

s, stato;

a, ingresso



## How About Learning the Value Function?



Facciamo imparare all'agente la value function, per una certa politica:  $V^\pi$ :

$$V^\pi(s) = \left[ \sum_{a_j} \pi(a_j, s) \right] \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^\pi(s')]$$

È una funzione dello stato.

Una volta imparata la value function,  $V^*$ , l'agente seleziona la policy ottima passo per passo, "one step lookahead":

$$\pi^*(s) = \arg \max_a \sum_{s'} P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V^*(s')]$$

*Full backup, for all states*



## Value function iteration



Facciamo imparare all'agente la value function, per una certa politica:  $V^\pi$ , analizzando quello che succede in uno step temporale:

$$V^{\pi}_{k+1}(s) = \left[ \sum_{a_j} \pi(a_j, s) \right] \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^{\pi}_k(s')]$$

L'apprendimento della policy si può inglobare nella value iteration:

$$V_{k+1}(s) = \max_a \sum_{s'} P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V_k(s')]$$

*Full backup, for all states*



## Asynchronous DP

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma V_k(s')]$$

*Full backup, single state, s, all future states s'*

Fino a questo punto, è noto un modello dell'ambiente:

- R(.)
- P(.)



## Temporal Differences

Viene rimossa ogni conoscenza a-priori sull'ambiente.

Value iteration and policy iteration are merged into one update equation:

$$V^\pi(s_t) = V^\pi(s_t) + \alpha [r_{t+1} + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)]$$

*Single backup, single state,  $s_t$ , single future state  $s_{t+1}$*

+ policy improvement



## Serve davvero la Value Function?



La Value Function deriva dalla visione della Programmazione Dinamica.

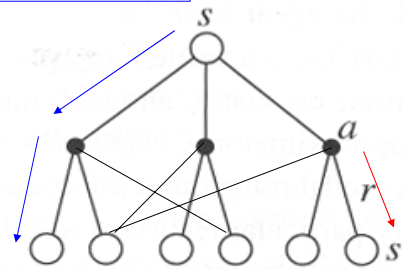
Ma è proprio necessario conoscere la Value function? In fondo a noi interessa determinare la Policy.



## Le value function



$$V^\pi(s) = E_\pi\{R_t | s_t = s\} = E_\pi\left\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s\right\} = \left[\sum_{a_j} \pi(a_j, s)\right] \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^\pi(s')]$$



$$Q^\pi(s, a) = E_\pi\{R_t | s_t = s, a_t = a\} = E_\pi\left\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a\right\}$$

$$= \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^\pi(s')]$$



## Equazioni di ottimalità di Bellman



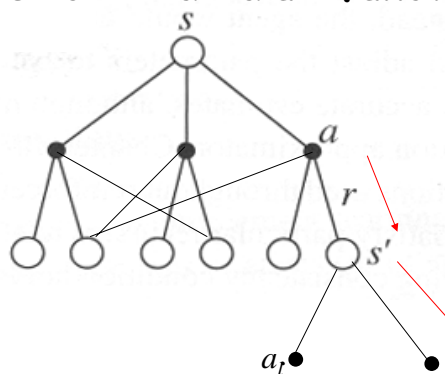
$V^*(s)$  di uno stato, quando viene scelta la policy ottima, deve essere uguale al valore atteso del reward per l'azione migliore per lo stato  $s$ .

$$V^*(s) = \max_{a_j} \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^*(s')] ]$$

$$Q^*(s, a_j) = \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma \max_{a'} Q^*(s', a')] ]$$



## Calcolo ricorsivo della value function $Q$



$$Q^\pi(s, a) = E_\pi \{ R_t | s_t = s, a_t = a \} = E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\}$$

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



## Q Functions

$$\pi^*(s) = \arg \max_a \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma V^\pi(s')] = \arg \max_a Q(s, a)$$

$V$  = Cumulative reward of being in  $s$  and choosing  $a_j$ ;  $Q^\pi(s, a_j) = \sum_{s'} P_{s \rightarrow s'}^{a_j} [R_{s \rightarrow s'}^{a_j} + \gamma V^\pi(s')]$

### Idea chiave:

- Unire il rinforzo che si ottiene passando da uno stato al successivo in un'unica funzione

$$Q(s, a) = [R_{s \rightarrow s'}^a + \gamma V^\pi(s')]$$

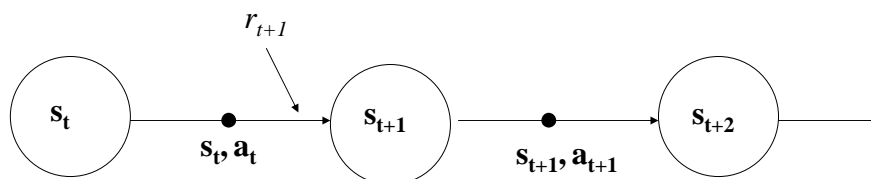
- Questa funzione valuta la bontà dell'azione e non più dello stato ( $a = \pi(s)$ ).
- A questo punto posso massimizzare  $Q$  senza conoscere separatamente il reward istantaneo e la value function come:

$$\pi^*(s) = \operatorname{argmax}_a Q(s, a)$$

$Q$  = Cumulative reward of being in  $s$  and taking action  $a$ .



## Rappresentazione grafica



$V(s_t)$

$V(s_{t+1})$

One step for Value Iteration

$Q(s_t, a_t)$

$Q(s_{t+1}, a_{t+1})$

One step for Q Iteration



## Sommario



Value e Q functions

**SARSA**

Q-learning



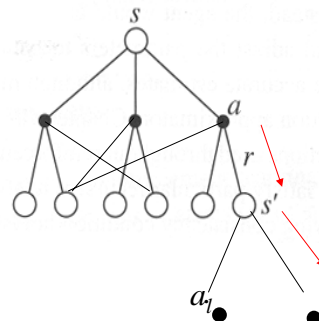
## Come apprendere Q: SARSA



$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

1) Apprendiamo il valore di Q per una policy data (on-policy).

2) Dopo avere appreso la funzione Q, possiamo modificare la policy in modo da migliorarla.







# SARSA Algorithm (progetto)



```

Q(s,a) = rand(); // ∀s, ∀a, eventualmente Q(s,a) = 0
Repeat // for each episode
{
  s = s0;
  Repeat // for each step of the single episode
  {
    a = Policy(s); // ε-greedy??
    s_next = NextState(s,a);
    reward = Reward(s,s_next,a);
    a_next = Policy(s_next); // ε-greedy?
    Q(s,a) = Q(s,a) + α [reward + γ Q(s_next, a_next) - Q(s,a)];
    s = s_next;
  } // until last state
} // until the end of learning

```

- 1) Apprendiamo il valore di Q per una policy data (on-policy).
- 2) Dopo avere appreso la funzione Q, possiamo modificare la policy in modo da migliorarla.

Come integrare i due passi?

A.A. 2006-2007

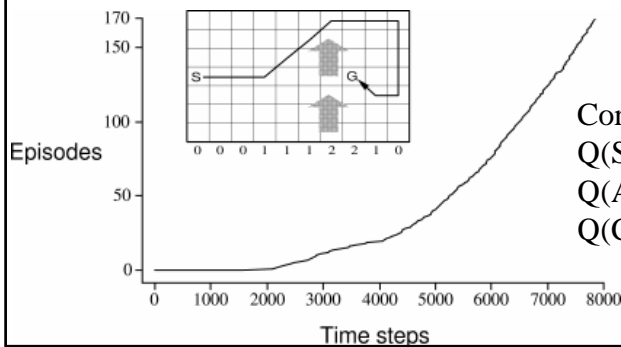
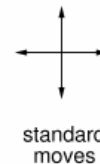
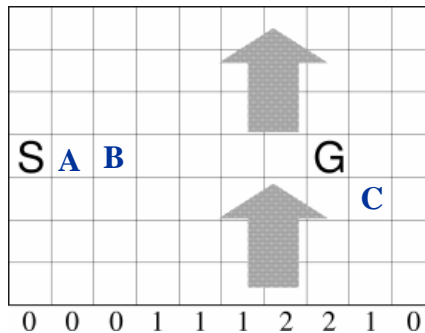
17/34

<http://homes.dsi.unimi.it/~borghese/>



## Esempio

$Q(s,a)$  iniziale = 0.01  
 $r = 0$  se  $s' = G$ ,  
 altrimenti  $r = -1$ .  
 $\pi(s,a)$  data.



Correzione di Q ad un passo:  
 $Q(S, \text{east}) = -0.09$   
 $Q(A, \text{east}) = -0.09$   
 $Q(C, \text{west}) = 0$

<http://homes.dsi.unimi.it/~borghese/>



## Sommario



Value e Q functions

SARSA

Q-learning



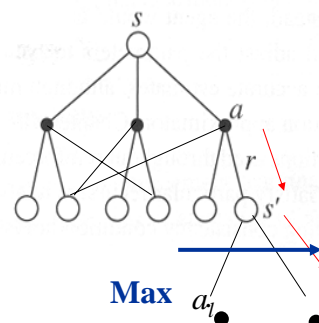
## Off-policy Temporal Difference: Q-learning



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

Non imparo semplicemente la funzione valore Q, ma la funzione valore Q ottima.

In  $s$ , scelgo un ramo del grafo, e poi **decido** ad un passo come continuare.





## Q-learning algorithm (progetto)



```

Q(s,a) = rand(); // ∀s, ∀a, Q(s,a) = 0 eventualmente
Repeat // for each episode
{
  s = s0;
  Repeat // for each step of the single episode
  {
    a = Policy(s); // eventualmente ε-greedy
    s_next = NextState(s,a);
    reward = Reward(s,s_next,a);
    a_next = Policy(s_next); // eventualmente ε-greedy
    Q(s,a) = Q(s,a) + α [reward + γ max Q(s_next, a_next) - Q(s,a)];
    s = s_next;
    UpdatePolicy(a_next, s_next);
  } // until last state
} // until the end of learning

```

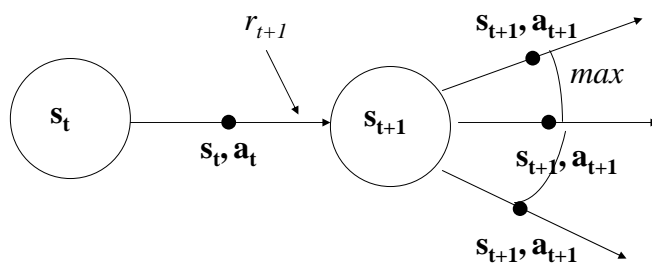
A.A. 2006-2007

21/34

<http://homes.dsi.unimi.it/~borghese/>



## Rappresentazione grafica



$Q(s_t, a_t)$   $Q(s_{t+1}, a_{t+1})$   
 One step for Q Iteration

Viene migliorata la policy al tempo  $t+1$ .

A.A. 2006-2007

22/34

<http://homes.dsi.unimi.it/~borghese/>



## Sommario

Dal modello dell'ambiente all'esperienza

**Q-learning**

Esempio

A.A. 2006-2007

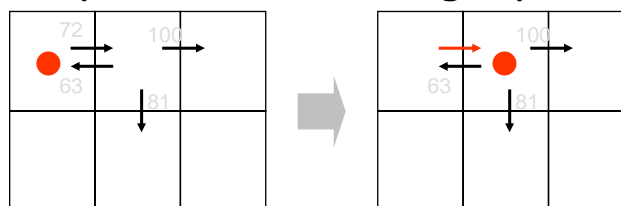
23/34

<http://homes.dsi.unimi.it/~borghese/>



## Example 1 - Q Learning Update

$\gamma = 0.9$



Esempio tratto dai lucidi del corso di Brian C. Williams su RL.

Modificati dalle slide di: Manuela Veloso, Reid Simmons, & Tom Mitchell, CMU

Apprendimento della funzione valore Q. Versione SARSA.

A.A. 2006-2007

24/34

<http://homes.dsi.unimi.it/~borghese/>

### Example 1 - Q Learning Update

$\gamma = 0.9$

$\alpha = 0.1$

0 reward received in the transition

$$\begin{aligned}
 Q(s_1, a_{right}) &\leftarrow Q(s_1, a_{right}) + \alpha \{ r(s_1, a_{right}, s_2) + \gamma \max_a Q(s_2, a) - Q(s_1, a_{right}) \} \\
 &\leftarrow 72 + \alpha [0 + 0.9 \max \{63, 81, 100\} - Q(s_1, a_{right}) ] \\
 &\leftarrow 72 + \alpha(90 - 72) = 1.8
 \end{aligned}$$

In grigio i valori di  $Q(s,a)$ . Nessun reward istantaneo.  
 Correzione di  $Q(s_1, a_{right})$   
 La correzione va a 0 quando  $Q(s_1, a_{right}) = 90$

A.A. 2006-2007 25/34 http://homes.dsi.unimi.it/~borghese/

### Example 2: Q-Learning Iterations: Episodic

- Start at upper left; Selected policy: move clockwise; Table initially 0;  $\gamma = 0.8$ .  
 Possibili transizione sono segnate con frecce nere e grigie.

Reward istanteo in rosso

$\alpha = 1$

$$Q(s_t, a_t) \leftarrow \left[ r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) \right]$$

Q(s1,E)	Q(s2,E)	Q(s3,S)	Q(s4,W)
0			

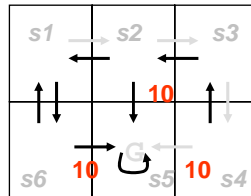
A.A. 2006-2007 26/34 http://homes.dsi.unimi.it/~borghese/



## Q-Learning Iterations

- Start at upper left – move clockwise; table initially 0;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$Q(s1,E)$	$Q(s2,E)$	$Q(s3,S)$	$Q(s4,W)$
0	0	0	

A.A. 2006-2007

27/34

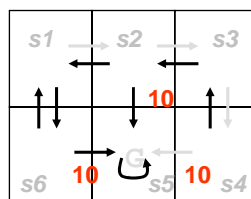
<http://homes.dsi.unimi.it/~borghese/>



## Q-Learning Iterations

- Start at upper left – move clockwise;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$Q(s1,E)$	$Q(s2,E)$	$Q(s3,S)$	$Q(s4,W)$
0	0	0	$r + \gamma \max_{a'} [Q(s5a)] = 10 + 0.8 \times 0 = 10$

A.A. 2006-2007

28/34

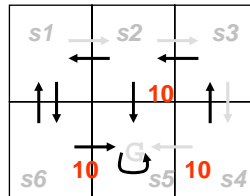
<http://homes.dsi.unimi.it/~borghese/>



## Q-Learning Iterations

- Start at upper left – move clockwise;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$Q(s1,E)$	$Q(s2,E)$	$Q(s3,S)$	$Q(s4,W)$
0	0	0	$r + \gamma \max_{a'} \{Q(s5,a)\} = 10 + 0.8 \times 0 = 10$
0	0	$r + \gamma \max_{a'} \{Q(s4,W), Q(s4,N)\} = 0 + 0.8 \times \max\{10,0\} = 8$	

A.A. 2006-2007

29/34

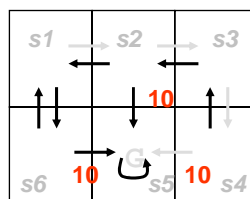
<http://homes.dsi.unimi.it/~borghese/>



## Q-Learning Iterations

- Start at upper left – move clockwise;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$Q(s1,E)$	$Q(s2,E)$	$Q(s3,S)$	$Q(s4,W)$
0	0	0	$r + \gamma \max_{a'} \{Q(s5,loop)\} = 10 + 0.8 \times 0 = 10$
0	0	$r + \gamma \max_{a'} \{Q(s4,W), Q(s4,N)\} = 0 + 0.8 \times \max\{10,0\} = 8$	10
0	$r + \gamma \max_{a'} \{Q(s3,W), Q(s3,S)\} = 0 + 0.8 \times \max\{0,8\} = 6.4$		

A.A. 2006-2007

30/34

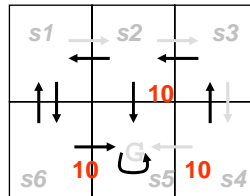
<http://homes.dsi.unimi.it/~borghese/>



## Q-Learning Iterations

- Start at upper left – move clockwise;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$Q(s1,E)$	$Q(s2,E)$	$Q(s3,S)$	$Q(s4,W)$
0	0	0	$r + \gamma \max_{a'} \{Q(s5,loop)\} = 10 + 0.8 \times 0 = 10$
0	0	$r + \gamma \max_{a'} \{Q(s4,W), Q(s4,N)\} = 0 + 0.8 \times \max\{10,0\} = 8$	10
0	$r + \gamma \max_{a'} \{Q(s3,W), Q(s3,S)\} = 0 + 0.8 \times \max\{0,8\} = 6.4$	8	10

A.A. 2006-2007

31/34

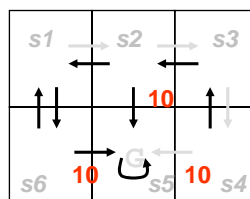
<http://homes.dsi.unimi.it/~borghese/>



## Q-Learning Iterations: improving policy

- Start at upper left – move clockwise;  $\gamma = 0.8$

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$



$$\text{Calcolo } Q(s_2, S) = r + \gamma \max_{a'} \{Q(s5,loop)\} = 10 + 0.8 \times 0 = 10$$


$$\text{Ricalcolo } Q(s1, E) = r + \gamma \max_{a'} \{Q(s2,E), Q(s2,W), Q(s2, S)\} = r + \gamma \max_{a'} \{6.4, 0.0, 10.0\} \rightarrow \text{South} = \pi(s_2)!$$

A.A. 2006-2007

32/34


<http://homes.dsi.unimi.it/~borghese/>

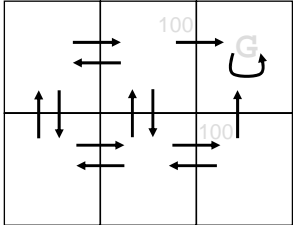




### Example 3

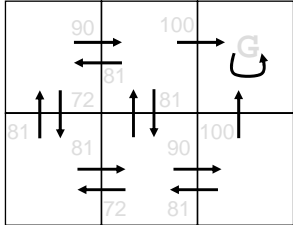
## Comparison functions V and Q



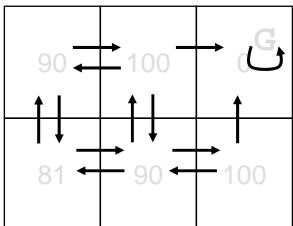


R(s, a) values

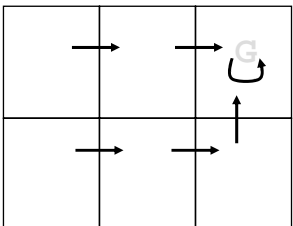
→



Q(s, a) values



V\*(s) values




One Optimal Policy


A.A. 2006-2007

33/34

<http://homes.dsi.unimi.it/~borghese/>



## Sommario



Value e Q functions

SARSA

Q-learning

A.A. 2006-2007

34/34

<http://homes.dsi.unimi.it/~borghese/>