

# Sistemi Intelligenti Reinforcement Learning: Temporal Differences

Alberto Borghese

Università degli Studi di Milano  
Laboratorio di Sistemi Intelligenti Applicati (AIS-La)  
Dipartimento di Informatica  
[alberto.borghese@unimi.it](mailto:alberto.borghese@unimi.it)  
*Barto and Sutton, Capitoli 3 e 6*



A.A. 2020-2021

1/48

<http://borghese.di.unimi.it/>



## Sommario



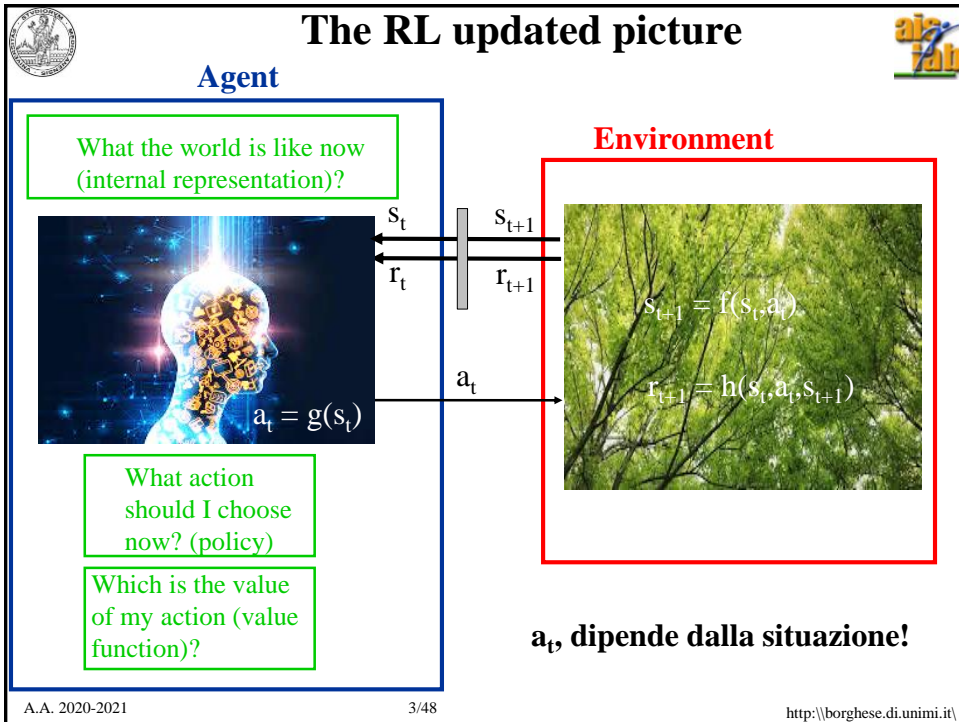
Le equazioni di Bellman

Differenze temporali

A.A. 2020-2021

2/48

<http://borghese.di.unimi.it/>



## Meccanismo di apprendimento nel RL




**Inizializzazione:** se l'agente non agisce sull'ambiente non succede nulla. Occorre specificare una policy iniziale.

**Ciclo dell'agente** (le tre fasi sono sequenziali):

- 1) Implemento una policy ( $\pi(s,a)$ )
- 2) Aggiorno la Value function ( $Q^\pi(s,a)$ )**
- 3) Aggiorno la policy.

A.A. 2020-2021

4/48

<http://borghese.di.unimi.it/>



## Esempio: AIBO search



### Azioni:

- 1) Rimanere fermo e aspettare che qualcuno getti nel cestino una lattina vuota.
- 2) Muoversi attivamente in cerca di lattine.
- 3) Tornare alla sua base (recharge station) e ricaricarsi.

### Stato:

- 1) Alto livello di energia.
- 2) Basso livello di energia.

**Goal:** collezionare il maggior numero di lattine.

### Azioni ammissibili (policy):

$a(s = \text{high}) = \{\text{Search, Wait}\}$

$a(s = \text{low}) = \{\text{Search, Wait, Recharge}\}$

A.A. 2020-2021

5/48

<http://borgnese.di.unimi.it/>

<http://borgnese.di.unimi.it/>



## Esempio di calcolo della Value function



Policy deterministica

$a(\text{high}) = \text{wait}$

$a(\text{low}) = \text{search}$

Value function

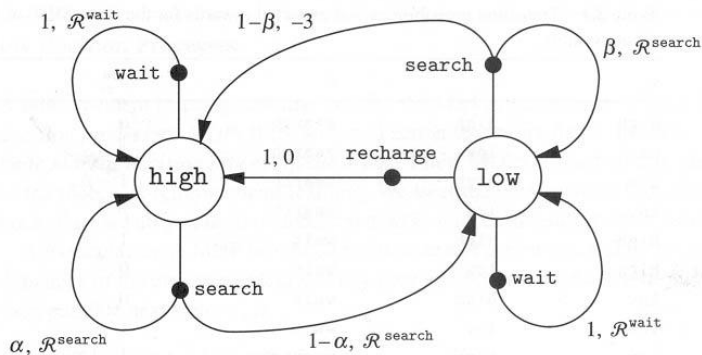
$Q(\text{high, search}) = ?$

$Q(\text{low, search}) = ?$

$\alpha = \Pr(s_{t+1} = \text{High} | s_t = \text{High}, a_t = \text{Search}) = 0.4$

$\beta = \Pr(s_{t+1} = \text{Low} | s_t = \text{Low}, a_t = \text{Search}) = 0.1$


$\gamma = 0.8, R^{\text{search}} = 3, R^{\text{wait}} = 1, R^{\text{dead}} = -3, R^{\text{auto}} = 0$




A.A. 2020-2021

6/48

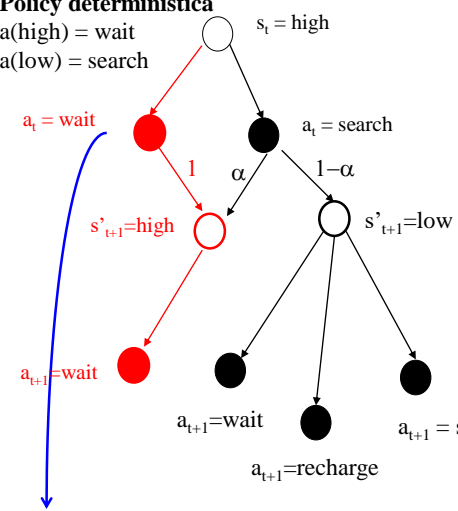
<http://borgnese.di.unimi.it/>



## Analisi ad un passo dal tempo t

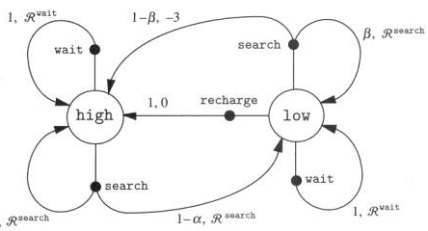


**Policy deterministica**  
 a(high) = wait  
 a(low) = search



$s_t = \text{high}$   
 $a_t = \text{wait}$   
 $s'_{t+1} = \text{high}$   
 $a_{t+1} = \text{wait}$

$a_t = \text{search}$   
 $s'_{t+1} = \text{low}$   
 $a_{t+1} = \text{wait}$   
 $a_{t+1} = \text{recharge}$   
 $a_{t+1} = \text{search}$



$1, R^{\text{wait}}$   
 $1-\beta, -3$   
 $\beta, R^{\text{search}}$   
 $1, 0$   
 $\alpha, R^{\text{search}}$   
 $1-\alpha, R^{\text{search}}$   
 $1, R^{\text{wait}}$


$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$$


A.A. 2020-2021

7/48

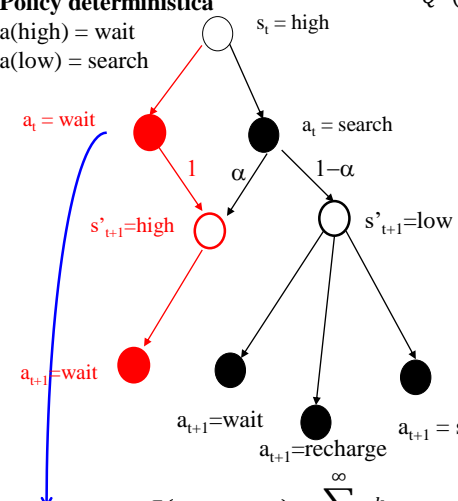
<http://borghese.di.unimi.it/>



## Analisi ad un passo dal tempo t



**Policy deterministica**  
 a(high) = wait  
 a(low) = search



$s_t = \text{high}$   
 $a_t = \text{wait}$   
 $s'_{t+1} = \text{high}$   
 $a_{t+1} = \text{wait}$

$a_t = \text{search}$   
 $s'_{t+1} = \text{low}$   
 $a_{t+1} = \text{wait}$   
 $a_{t+1} = \text{recharge}$   
 $a_{t+1} = \text{search}$

$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$$

$$Q^\pi(\text{high}, \text{wait}) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} =$$

$$R^{\text{wait}} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} =$$

$$R^{\text{wait}} + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k+1} =$$

$$R^{\text{wait}} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2}$$


$$Q^\pi(\text{high}, \text{wait}) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} = R^{\text{wait}} + \gamma Q^\pi(\text{high}, \text{wait})$$

$$Q^\pi(h, w) = [1 + 0.8 Q^\pi(h, w)]$$


A.A. 2020-2021

8/48

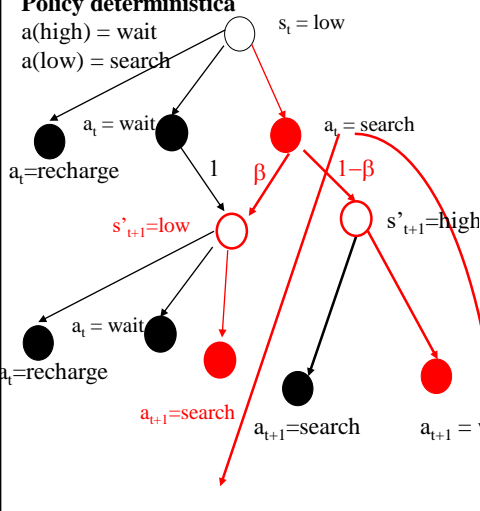
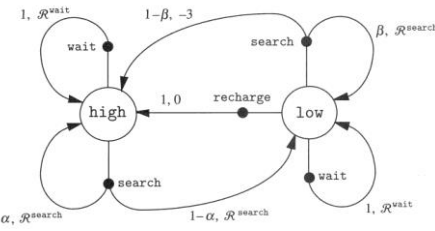
<http://borghese.di.unimi.it/>



## Analisi ad un passo dal tempo t



**Policy deterministica**  
 $a(\text{high}) = \text{wait}$   
 $a(\text{low}) = \text{search}$





**2 cammini possibili!!**


$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$$

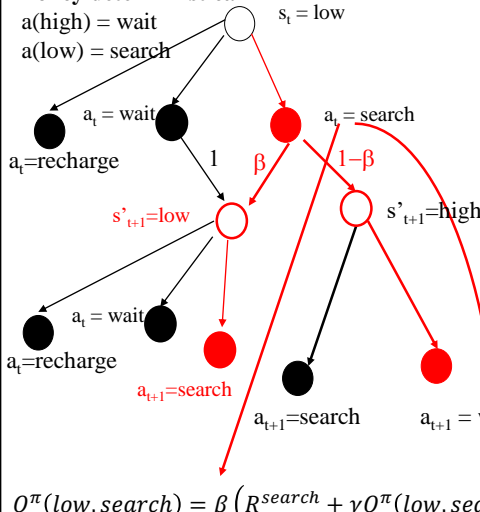
A.A. 2020-2021 9/48 http://borghese.di.unimi.it/



## Analisi ad un passo dal tempo t



**Policy deterministica**  
 $a(\text{high}) = \text{wait}$   
 $a(\text{low}) = \text{search}$



**2 cammini possibili!!**

- 1)  $R^{\text{search}} + \gamma Q^\pi(\text{low}, \text{search})$
- 2)  $R^{\text{dead}} + \gamma Q^\pi(\text{low}, \text{high})$

$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(\text{low}, \text{search}) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} =$$

$$Q^\pi(\text{low}, \text{search}) = \beta (R^{\text{search}} + \gamma Q^\pi(\text{low}, \text{search})) + (1 - \beta) (R^{\text{dead}} + \gamma Q^\pi(\text{high}, \text{wait}))$$

$$Q(1, s) = 0.1 \times [3 + 0.8 \times Q(1, s)] + 0.9 \times [-3 + 0.8 \times Q(h, w)]$$

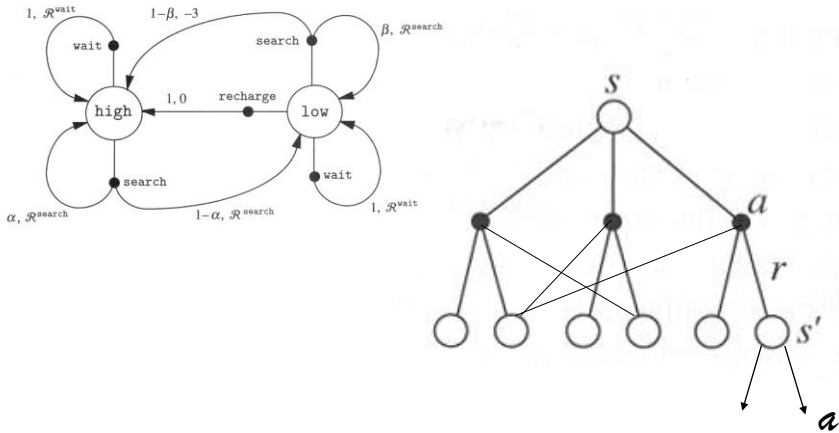
Contiene la probabilità di ricevere un reward  $\gamma Q(s', a)$ , condizionata a  $s_{t+1} = s'$  ip://borghese.di.unimi.it/



# Valutazione policy stocastica



Nel valutare  $Q(s,a)$  dobbiamo valutare tutti i cammini che partono da ogni  $s'$ .



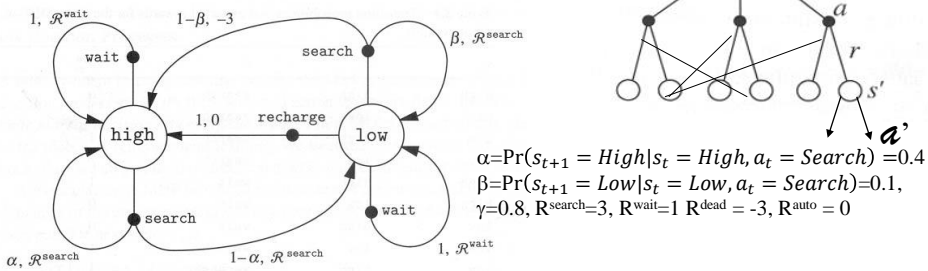
A.A. 2020-2021

11/48

<http://borghese.di.unimi.it>



# Policy stocastica



$$Q(\text{high}, \text{wait}) = 1 \times \{R_{wait} + \gamma [\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait})]\}$$

$$Q(\text{high}, \text{wait}) = 1 \times \{1 + 0.8 [\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait})]\}$$

$$Q(\text{high}, \text{search}) = \Pr(s_{t+1} = High | s_t = High, a_t = Search) \times$$

$$\{R_{search} + \gamma [\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait})]\} +$$

$$(1 - \Pr(s_{t+1} = High | s_t = High, a_t = Search)) \times$$

$$\{R_{search} + \gamma [\Pr(a' = \text{search} | \text{low}) Q(\text{low}, \text{search}) + \Pr(a' = \text{wait} | \text{low}) Q(\text{low}, \text{wait}) + \Pr(a' = \text{recharge} | \text{low}) Q(\text{low}, \text{rech})]\}$$

$$Q(\text{high}, \text{search}) = 0.4 \times \{3 + 0.8 [\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait})]\} +$$

$$0.6 \times \{3 + 0.8 [\Pr(a' = \text{search} | \text{low}) Q(\text{low}, \text{search}) + \Pr(a' = \text{wait} | \text{low}) Q(\text{low}, \text{wait})$$

$$+ \Pr(a' = \text{recharge} | \text{low}) Q(\text{low}, \text{rech})]\}$$

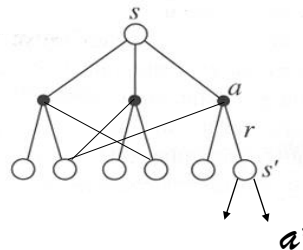
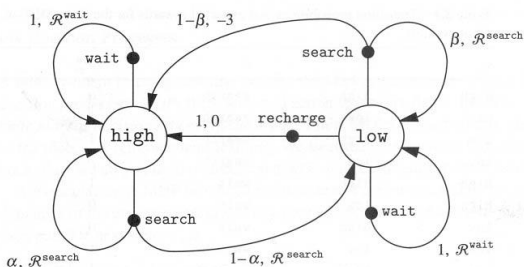
A.A. 2020-2021

12/48

<http://borghese.di.unimi.it>



# Policy stocastica



$\alpha=0.4, \beta=0.1, \gamma=0.8,$   
 $\mathcal{R}^{search}=3, \mathcal{R}^{wait}=1, \mathcal{R}^{dead}=-3, \mathcal{R}^{auto}=0$

$$Q(\text{low}, \text{wait}) = 1 \times \{ \mathcal{R}^{wait} + \gamma [ \text{Pr}(a'=\text{search}) Q(\text{low}, \text{search}) + \text{Pr}(a'=\text{wait}) Q(\text{low}, \text{wait}) + \text{Pr}(a'=\text{recharge}) Q(\text{low}, \text{recharge}) ] \}$$

$$Q(\text{low}, \text{search}) = \beta \times \{ \mathcal{R}^{search} + \gamma [ (\text{Pr}(a'=\text{search}) Q(\text{high}, \text{search}) + \text{Pr}(a'=\text{wait}) Q(\text{high}, \text{wait}) + \text{Pr}(a'=\text{recharge}) Q(\text{low}, \text{recharge})) ] \} + (1-\beta) \times \{ \mathcal{R}^{dead} + \gamma [ \text{Pr}(a'=\text{search}) Q(\text{high}, \text{search}) + \text{Pr}(a'=\text{wait}) Q(\text{high}, \text{wait}) ] \}$$

$$Q(\text{low}, \text{recharge}) = 1 \times \{ \mathcal{R}^{auto} + \gamma [ (\text{Pr}(a'=\text{search}) Q(\text{high}, \text{search}) + \text{Pr}(a'=\text{wait}) Q(\text{high}, \text{wait}) ) ] \}$$

A.A. 2020-2021

5 equazioni in 5 incognite

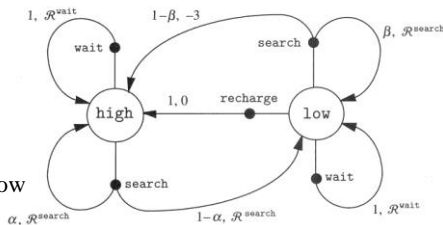
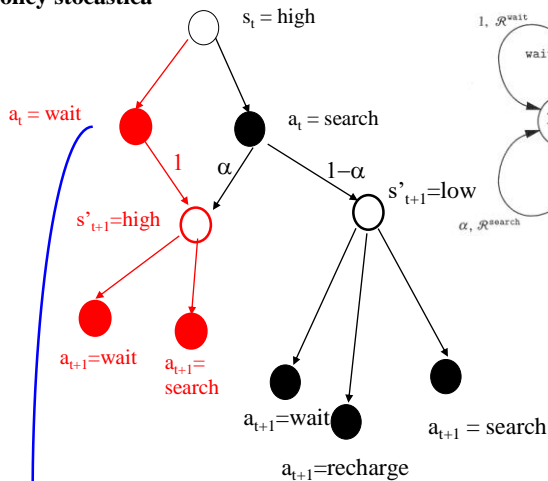
<http://borghese.di.unimi.it/>



# Analisi ad un passo dal tempo t



Policy stocastica



$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(s_t, a_t) = E_\pi \{ R_t | s_t = s, a_t = a \}$$

A.A. 2020-2021

14/48

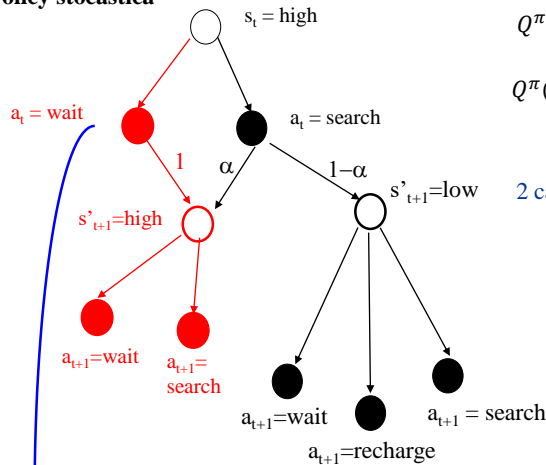
<http://borghese.di.unimi.it/>



# Analisi ad un passo dal tempo t



Policy stocastica



$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(s_t, a_t) = E_{\pi}\{R_t | s_t = s, a_t = a\}$$

2 cammini possibili!!

$$1) R^{wait} + \gamma Q^\pi(high, wait)$$

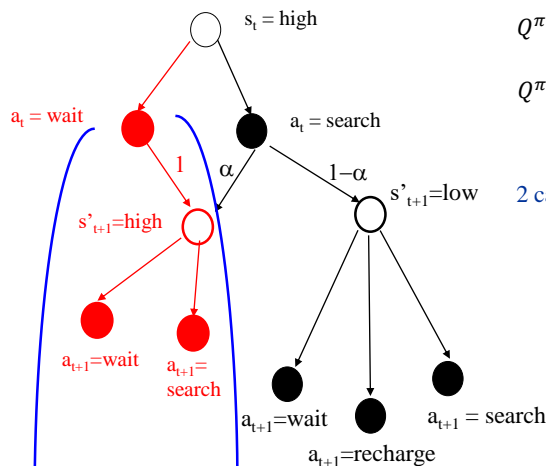
$$2) R^{wait} + \gamma Q^\pi(high, search)$$



# Analisi ad un passo dal tempo t



Policy stocastica (uniforme)



$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$Q^\pi(s_t, a_t) = E_{\pi}\{R_t | s_t = s, a_t = a\}$$


2 cammini possibili!!

$$1) R^{wait} + \gamma Q^\pi(high, wait)$$


$$2) R^{wait} + \gamma Q^\pi(high, search)$$

$$Q^\pi(high, wait) = R^{wait} + 0.5 \gamma Q^\pi(high, wait) + 0.5 \gamma Q^\pi(high, search)$$

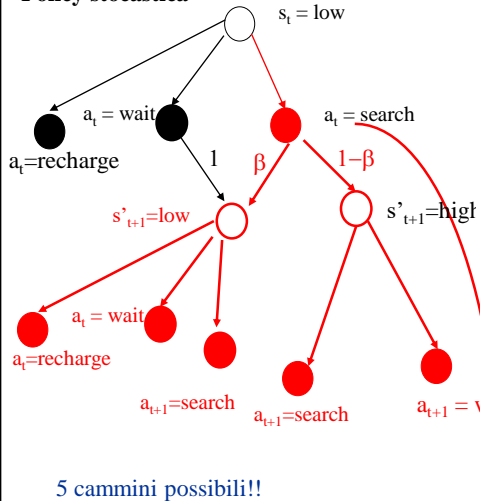




## Analisi ad un passo dal tempo t



**Policy stocastica**



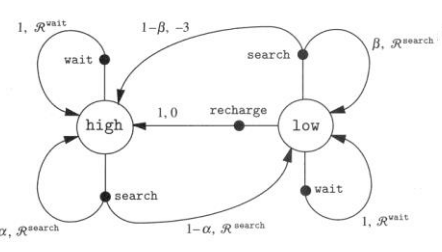
$s_t = \text{low}$

$a_t = \text{wait}$ ,  $a_t = \text{recharge}$ ,  $a_t = \text{search}$

$s'_{t+1} = \text{low}$ ,  $s'_{t+1} = \text{high}$

$a_{t+1} = \text{wait}$ ,  $a_{t+1} = \text{recharge}$ ,  $a_{t+1} = \text{search}$

5 cammini possibili!!



$1, R^{\text{wait}}$ ,  $1-\beta, -3$ ,  $\beta, R^{\text{search}}$

$1, 0$ ,  $1-\alpha, R^{\text{search}}$ ,  $1, R^{\text{wait}}$ ,  $\alpha, R^{\text{search}}$

$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$


$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$$

$$Q^\pi(\text{low}, \text{search}) = E_\pi\{R_t | s_t = \text{low}, a_t = \text{search}\}$$


A.A. 2020-2021

17/48

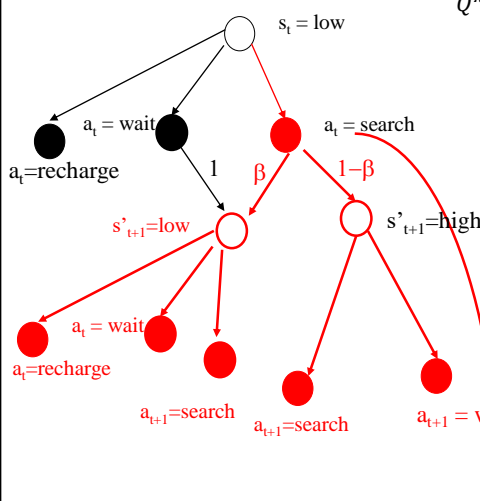
<http://borghese.di.unimi.it/>



## Analisi ad un passo dal tempo t



**Policy stocastica (equiprobabile)**



$s_t = \text{low}$

$a_t = \text{wait}$ ,  $a_t = \text{recharge}$ ,  $a_t = \text{search}$

$s'_{t+1} = \text{low}$ ,  $s'_{t+1} = \text{high}$

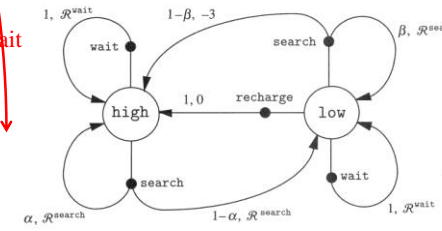
$a_{t+1} = \text{wait}$ ,  $a_{t+1} = \text{recharge}$ ,  $a_{t+1} = \text{search}$

5 cammini possibili!!

$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$$

A)  $R^{\text{search}} + \gamma[\frac{1}{3}Q^\pi(\text{low}, \text{search}) + \frac{1}{3}Q^\pi(\text{low}, \text{wait}) + \frac{1}{3}Q^\pi(\text{low}, \text{recharge})]$

B)  $R^{\text{dead}} + \gamma[\frac{1}{2}Q^\pi(\text{high}, \text{search}) + \frac{1}{2}Q^\pi(\text{high}, \text{wait})]$




$1, R^{\text{wait}}$ ,  $1-\beta, -3$ ,  $\beta, R^{\text{search}}$

$1, 0$ ,  $1-\alpha, R^{\text{search}}$ ,  $1, R^{\text{wait}}$ ,  $\alpha, R^{\text{search}}$

A.A. 2020-2021


18/48

<http://borghese.di.unimi.it/>



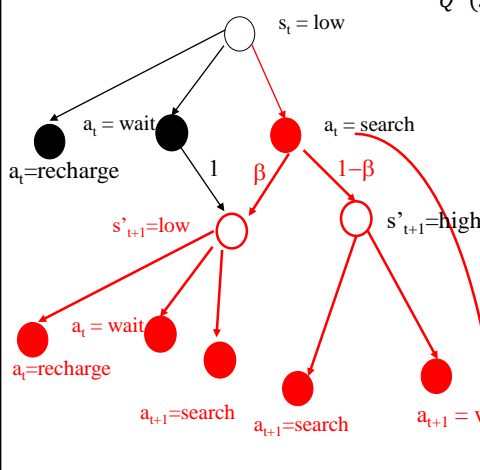
## Analisi ad un passo dal tempo t

Policy stocastica (equiprobabile)



$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\}$

5 cammini possibili!!




A)  $R^{search} + \gamma[\frac{1}{3}Q^\pi(low, search) + \frac{1}{3}Q^\pi(low, wait) + \frac{1}{3}Q^\pi(low, recharge)]$

B)  $R^{dead} + \gamma[\frac{1}{2}Q^\pi(high, search) + \frac{1}{2}Q^\pi(high, wait)]$


$Q^\pi(low, search) = \beta[R^{search} + \gamma(\frac{1}{3}Q^\pi(low, search) + \frac{1}{3}Q^\pi(low, wait) + \frac{1}{3}Q^\pi(low, recharge))] + (1-\beta)[R^{dead} + \gamma(\frac{1}{2}Q^\pi(high, search) + \frac{1}{2}Q^\pi(high, wait))]$

5 equazioni in 5 incognite

A.A. 2020-2021
19/48
<http://borghese.di.unimi.it/>



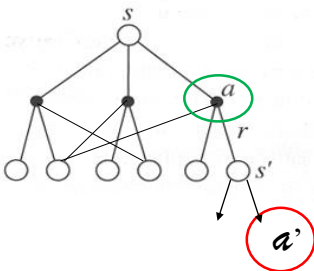
## Calcolo ricorsivo della Value function



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

$$Q^\pi(s_{t+1}, a_{t+1}) = E_\pi\{R_t | s_{t+1} = s', a_{t+1} = a'\}$$

Relazione tra  $Q^\pi(s, a)$  e  $Q^\pi(s', a')$ ?



A.A. 2020-2021
20/48
<http://borghese.di.unimi.it/>



# Calcolo ricorsivo della Value function



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

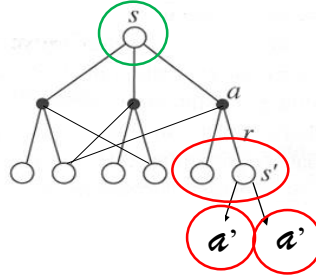
Isolo il reward ad un passo nella serie dei reward.

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

$$Q^\pi(s_t, a_t) = E_\pi\left\{\gamma^0 r_{t+1} + \sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} | s_t = s, a_t = a\right\}$$

Io termine  
(a un passo)

Io termine  
(passi futuri)



# $Q^\pi(s, a)$ : primo termine

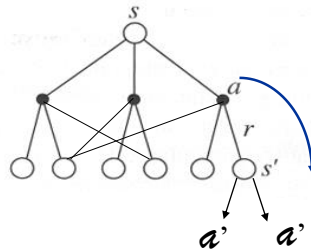


$$P_{s \rightarrow s' | a} \triangleq \Pr(s_{t+1} = s' | s_t = s, a_t = a)$$

$$E_\pi\{r_{t+1} | s_t = s, a_t = a\} = \sum_{s'} P_{s \rightarrow s' | a} R_{s, s', a}$$

Per ogni stato-azione devo valutare:

- Più stati prossimi
- Reward stocastici nella transizione ad un passo



**Visione Statistica:** Probabilità di ottenere il reward:  
condizionata all'arrivare nello stato  $s'$ :  $R_{s \rightarrow s' | a_j}$

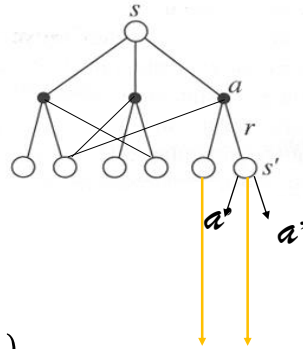


## $Q^\pi(s, a)$ : secondo termine



$$E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s_t = s, a_t = a \right\}$$

$$P_{s \rightarrow s' | a} \triangleq \Pr(s_{t+1} = s' \mid s_t = s, a_t = a)$$



$$E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^{k+1} r_{t+k+2} \mid s_t = s, a_t = a \right\}$$

$$= \gamma \sum_{s'} P_{s \rightarrow s' | a} E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_{t+1} = s' \right\}$$



## Putting all together



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

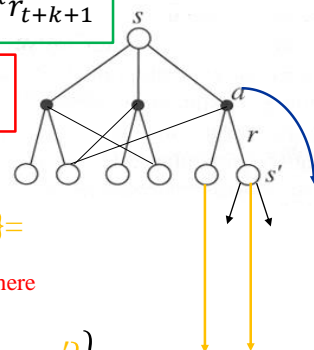
$$Q^\pi(s_{t+1}, a_{t+1}) = E_\pi \{ R_t \mid s_{t+1} = s', a_{t+1} = a' \}$$

$$\sum_{s'} P_{s \rightarrow s' | a} R_{s, s', a} + \gamma \sum_{s'} P_{s \rightarrow s' | a} E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_{t+1} = s' \right\} =$$

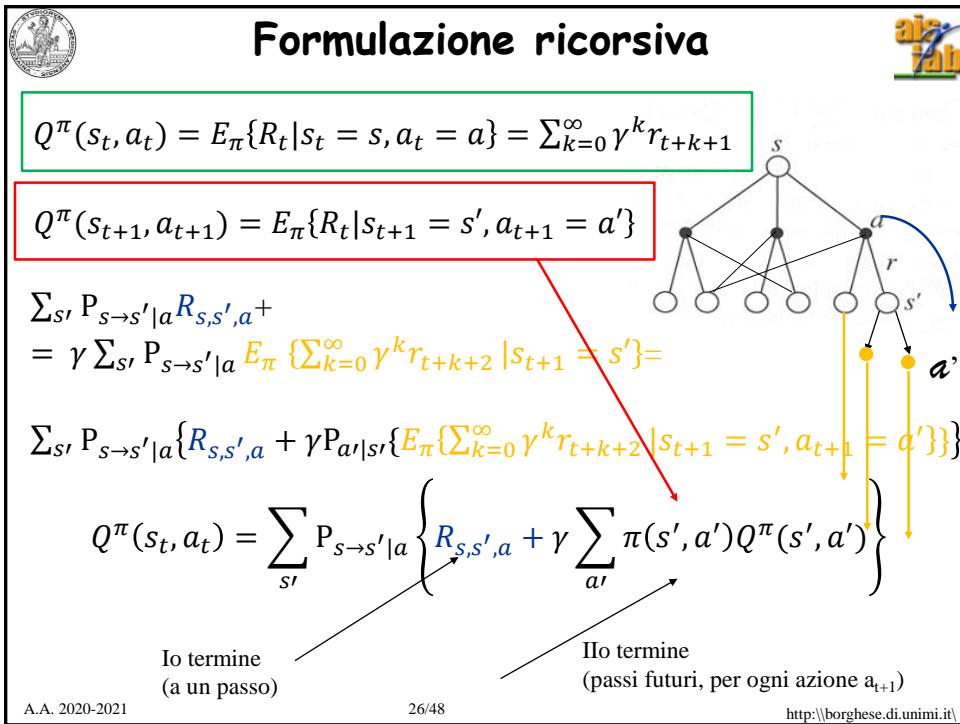
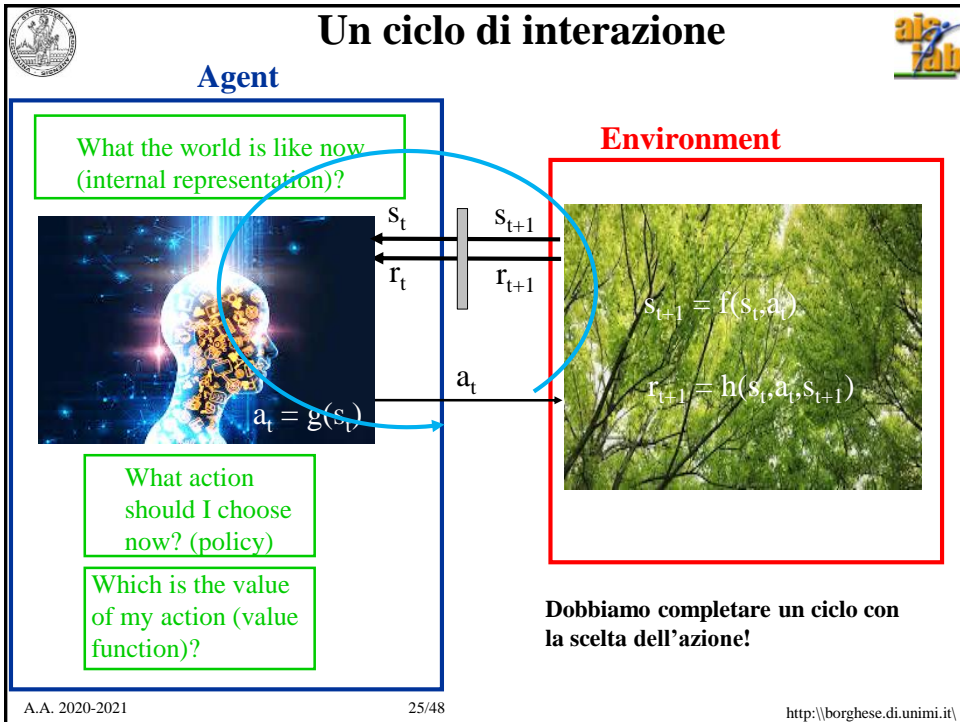
$$\sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_{t+1} = s' \right\} \right\}$$

Io termine  
(a un passo)

Il termine  
(passi futuri)



Not yet there

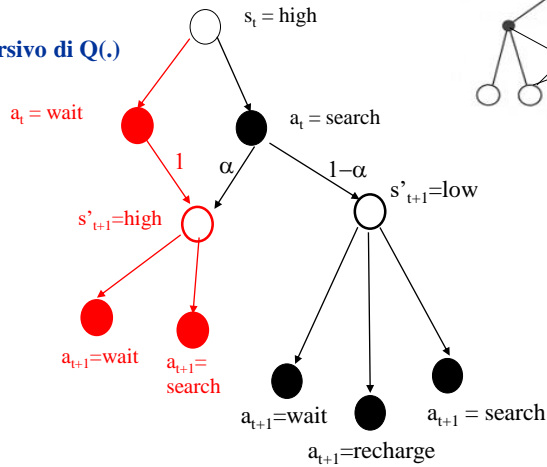




# Equazioni di Bellman

$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s \rightarrow s' | a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Calcolo ricorsivo di Q(.)



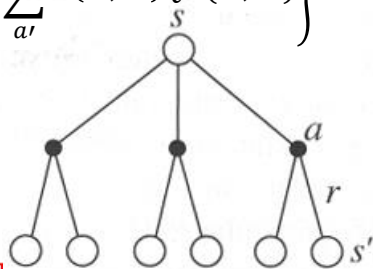
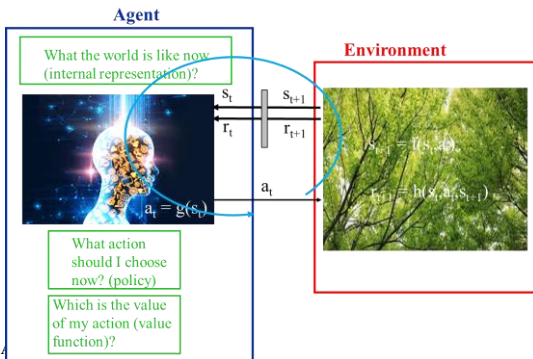
$$Q^\pi(\text{high}, \text{wait}) = R^{\text{wait}} + 0.5 \gamma Q^\pi(\text{high}, \text{wait}) + 0.5 \gamma Q^\pi(\text{high}, \text{search})$$



# Osservazioni

$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Calcolo ricorsivo di Q(.)



Passo da t a t+1 poi guardo backwards in time

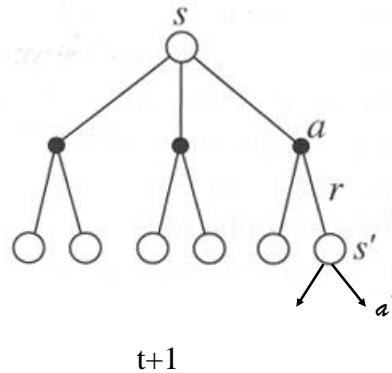
<http://borgnese.di.unimi.it/>



## Tecnica full-back

Back-up  
↑

$\pi(s,a)$  fissata



Conosciamo  $Q(s_t, a_t) \forall s_t, a_t$  anche per  $\{s'_{t+1}, a'_{t+1}\}$  quindi:

- Analizziamo la transizione da  $\{s_t, a_t\} \rightarrow \{s'_{t+1}, a'_{t+1}\}$
- Calcoliamo un nuovo valore di  $Q$  per  $\{s, a: Q(s_t, a_t)$  congruente con:

$Q(s_t, a_t)$  ed  $r_{t+1}$

*Full backup* se esaminiamo tutti gli  $s'$  e  $a'$  (cf. DP).

Da  $\{s', a'\}$  mi guardo indietro e aggiorno  $Q(s, a)$ .

$\pi$  fissata



## Meccanismo di apprendimento nel RL



**Inizializzazione:** se l'agente non agisce sull'ambiente non succede nulla. Occorre specificare una policy iniziale.

**Ciclo dell'agente (le tre fasi sono sequenziali):**

- 1) Implemento una policy ( $\pi(s,a)$ )
- 2) Aggiorno la Value function ( $Q^\pi(s,a)$ )
- 3) **Aggiorno la policy.**



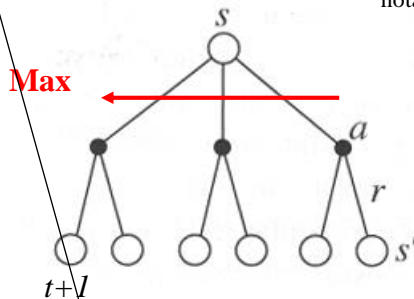
## Q(s, a) - Osservazioni

$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Policy nota

Per ogni stato devo valutare con informazioni esclusivamente racchiuse in 1 passo l'azione migliore a lungo termine

$$a_{new} : \max_a Q(s, a)$$



E' supposto noto il funzionamento dell'ambiente (simulazione)



## Sommario

Le equazioni di Bellman

Differenze temporali





## Q(s, a) - Osservazioni

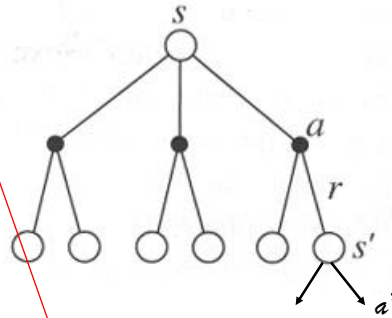


$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Policy nota

Per ogni stato devo valutare con informazioni esclusivamente racchiuse in 1 passo l'azione migliore a lungo termine

$$a_{new} : \max_a Q(s, a)$$



Non è noto il funzionamento dell'ambiente (interazione)



## Background su Temporal Difference (TD) Learning



Al tempo  $t$  abbiamo a disposizione:

$r_{t+1} = r'$  estratto (sampled) dalla distribuzione statistica:  $R_{s \rightarrow s' | a_j}$

$s_{t+1} = s'$  estratto (sampled) dalla distribuzione statistica:  $P_{s \rightarrow s' | a_j}$

**Dopo la realizzazione di un evento, l'incertezza statistica scompare.**

- 1 Reward certo
  - 1 Transizione certa
- vengono forniti dall'ambiente

Come si possono utilizzare per apprendere?



## Confronto con il rinforzo classico



$$Q_{k+1} = Q_k - \frac{Q_k}{N_{k+1}} + \frac{r_{k+1}}{N_{k+1}} = Q_k + \alpha[r_{k+1} - Q_k]$$

Occupazione di memoria minima: Solo  $Q_k$  e  $k$ .  
NB  $N_k$  è il numero di volte in cui è stata scelta  $a_j$ .

Questa forma è la base del RL. La sua forma generale è:

$$\begin{aligned} \text{NewEstimate} &= \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}] \\ \text{NewEstimate} &= \text{OldEstimate} + \text{StepSize} * \text{Error} \end{aligned}$$

$$\text{StepSize} = \alpha = 1/(N+1)$$

$$a = \text{cost}$$

$$\text{Rewards weight } w = 1$$

$$\text{Weight of } i\text{-th reward at time } k: w = (1-\alpha)^{k-i}$$

Qual è la differenza introdotta dall'approccio che prevede comportamenti (catene di azioni)?

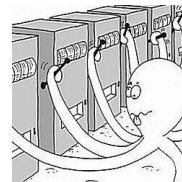


## Un possibile aggiornamento di $Q(s, a)$



$$Q_{k+1}(a) = Q_k(a) - \frac{Q_k(a)}{N_{k+1}(a)} + \frac{r_{k+1}(a)}{N_{k+1}(a)} = Q_k(a) + \alpha[r_{k+1}(a) - Q_k(a)] =$$

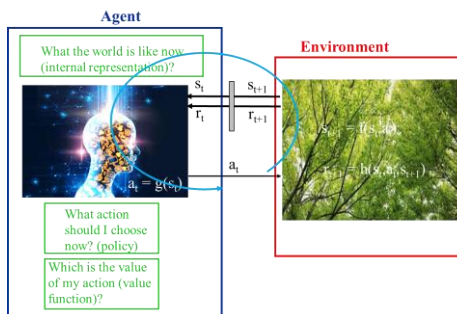
$$Q_k(a) + \alpha \Delta Q_k(a)$$




Come passo ai comportamenti?


$$Q_{k+1}^\pi(s, a) = Q_k^\pi(s, a) + \alpha \Delta Q_k(s, a)$$

Come calcolo  $\Delta Q_k$ ?





## Calcolo di $\Delta Q_k$



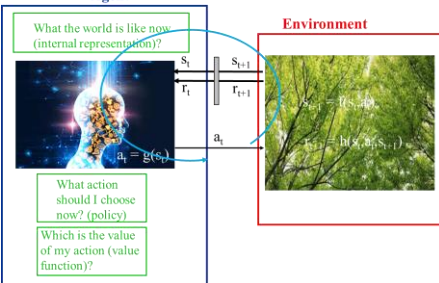
$$Q_{k+1}(a) = Q_k(a) + \alpha[r_{k+1}(a) - Q_k(a)] =$$

$$Q_k(a) + \alpha \Delta Q_k(a)$$

Al tempo  $t$  abbiamo a disposizione:

$r_{t+1} = r'$  da:  $R_{s \rightarrow s' | a_j}$

$s_{t+1} = s'$  da:  $P_{s \rightarrow s' | a_j}$



Quale semantica hanno  $Q(s,a)$  e  $r(s,a,s')$  nel caso dei comportamenti?

$$Q_{k+1}^\pi(s, a) = Q_k^\pi(s, a) + \alpha[r' + \gamma Q_k^\pi(s', a') - Q_k^\pi(s, a)] =$$

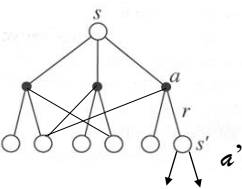
$$Q_k(a) + \alpha \Delta Q_k(a)$$

↑

Reward a 1 passo

↑


Reward a lungo termine da  $s'$




<http://borgnese.di.unimi.it/>

A.A. 2020-2021

37/48

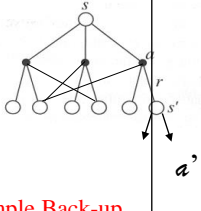


## TD(0) update



Ad ogni istante di tempo di ogni trial aggiorniamo la Value function:

$$Q_{k+1}^\pi(s, a) = Q_k^\pi(s, a) + \alpha[r' + \gamma Q_k^\pi(s', a') - Q_k^\pi(s, a)]$$



Sample Back-up

Conosciamo  $Q(s_t, a_t) \forall s_t, a_t$  anche per  $\{s'_{t+1}, a'_{t+1}\}$  quindi:

- Analizziamo la transizione da  $\{s_t, a_t\} \rightarrow \{s'_{t+1}, a'_{t+1}\}$
- Calcoliamo un nuovo valore di  $Q$  per  $\{s, a\}$  congruente con:
 
$$Q(s_t, a_t) \text{ ed } r_{t+1}$$

*Sample backup* se esaminiamo una sola coppia di  $s'$  e  $a'$  (cf. DP asincrona).  
Da  $\{s', a'\}$  mi guardo indietro e aggiorno  $Q(s, a)$ .

**Per  $\alpha$  che diminuisce con l'apprendimento, per  $k \rightarrow \infty$ ,  $Q_k^\pi(s, a)$  converge al valore vero di  $Q^\pi(s, a)$**

$\pi(s, a)$  fissata

**Posso ragionare a un passo per calcolare  $Q^\pi(s, a)$**

A.A. 2020-2021

38/48



## Confronto con il setting associativo



$$Q_{k+1} = Q_k - \frac{Q_k}{N_{k+1}} + \frac{r_{k+1}}{N_{k+1}} = Q_k + \alpha [r_{k+1} - Q_k]$$

Occupazione di memoria minima: Solo  $Q_k$  e  $k$ .  
NB  $k$  è il numero di volte in cui è stata scelta  $a_j$ .

Questa forma è la base del RL. La sua forma generale è:

$$\text{NewEstimate} = \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

$$\text{NewEstimate} = \text{OldEstimate} + \text{StepSize} * \text{Error}.$$

$$\text{StepSize} = \alpha = 1/N_{k+1} \quad a = \text{cost}$$



## Setting $\alpha$ value



$\alpha(s_t, a_t, s_{t+1}) = \frac{1}{N(s_t, a_t, s_{t+1})}$ , where  $N(s_t, a_t, s_{t+1})$  represents the number of occurrences of  $s_t, a_t, s_{t+1}$ . With this setting the estimated  $Q$  tends to the expected value of  $Q(s, a)$ .

**Per semplicità si assume solitamente  $\alpha < 1$  costante.** In questo caso,  $Q(s, a)$  assume il valore di una media pesata dei reward a lungo termine collezionati a partire da  $(s, a)$ , con peso:  $(1-\alpha)^k$ : *exponential recency-weighted average*.

**$\alpha$  che decresce dolcemente a zero consente la convergenza del Sistema stocastico.**



## Esempio



Stima del tempo di percorrenza da casa all'ufficio su un percorso ben definite (policy determinate e deterministica).

La durata dei diversi segmenti può variare da giorno a giorno e quindi la stima della durata totale viene corretta conseguentemente.

La stima corrente del tempo totale è data dalla somma dei tempi per:

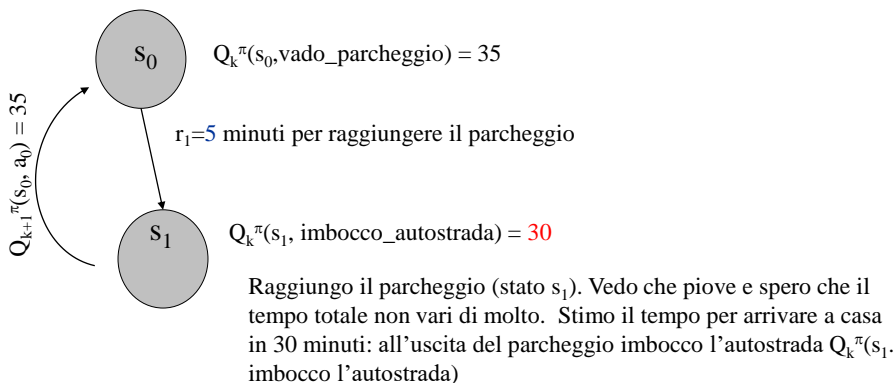
- Dall'ufficio al parcheggio: 5 minuti (time to go = 35 minuti)
  - Dal parcheggio all'uscita dell'autostrada: 15 minuti (time to go = 30 minuti)
  - Dall'uscita dell'autostrada alla strada di casa: 5 minuti (time to go = 15 minuti)
  - Dalla strada di casa a casa: 7 minuti (time to go = 10 minuti)
  - Dal parcheggio a casa: 3 minuti (time to go = 3 minuti)
- In totale 35 minuti.



## Learning $Q^\pi(s, a)$ - I



$s_0$  = ufficio;  $Q_k^\pi(s_0, \text{vado\_parcheggio}) = 35 \text{ minuti}$ ;  $Q_{k+1}^\pi(s_0, \text{vado\_parcheggio}) = 35 \text{ minuti}$  (potrei fare altre scelte, e.g. andare alla metropolitana, ma la policy prescrive di andare a prendere l'auto nel parcheggio perchè era considerata la soluzione più veloce).



Aggiorno il tempo totale, ovvero il tempo dallo stato  $s_0$ :

$$Q_{k+1}^\pi(s_0, a) = Q_k^\pi(s_0, a) + \alpha[r' + \gamma Q_k^\pi(s_1, a') - Q_k^\pi(s_0, a)] = 35 + [5 + 30 - 35] = 35$$

Suppongo  $\alpha = \gamma = 1$



## Learning $Q^\pi(s, a)$ - II

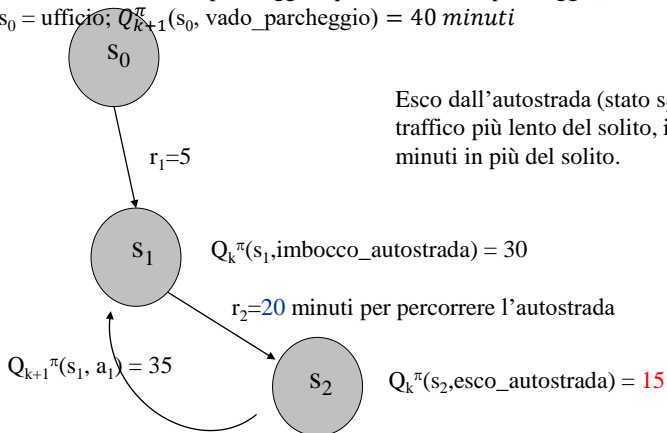


$s_1 = \text{parcheggio}$ ;  $Q_k^\pi(s_1, \text{imbocco\_autostrada}) = 30 \text{ minuti}$ ;

$Q_{k+1}^\pi(s_1, \text{imbocco\_autostrada}) = 30 \text{ minuti}$ ;

(potrei fare altre scelte, e.g. tornare in ufficio; una volta scelto di uscire, aggiorno il valore dell'azione uscire dal parcheggio, quando sono nel parcheggio)

$s_0 = \text{ufficio}$ ;  $Q_{k+1}^\pi(s_0, \text{vado\_parcheggio}) = 40 \text{ minuti}$



Aggiorno il tempo totale dallo stato  $s_1$ :

$$\Delta. Q_{k+1}^\pi(s_1, a) = Q_k^\pi(s_1, a) + \alpha[r' + \gamma Q_k^\pi(s_2, a') - Q_k^\pi(s_1, a)] = 30 + [20 + 15 - 30] = 35$$

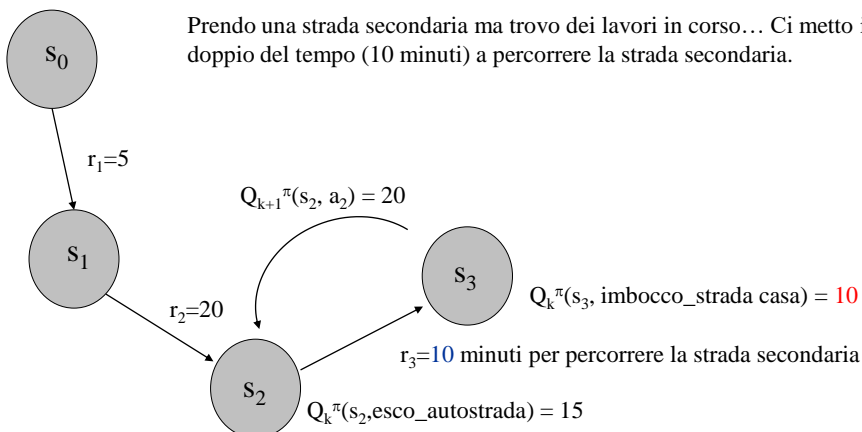


## Learning $Q^\pi(s, a)$ - III



$s_2 = \text{esco\_autostrada}$ ;  $Q_k^\pi(s_2, \text{esco\_autostrada}) = 15 \text{ min}$ ;  $Q_{k+1}^\pi(s_2, \text{esco\_autostrada}) = 20 \text{ min}$ ;

$s_0 = \text{ufficio}$ ;  $Q_{k+1}^\pi(s_0, \text{vado\_parcheggio}) = 45 \text{ minuti}$



Aggiorno il tempo totale dallo stato  $s_2$ :

$$\Delta. Q_{k+1}^\pi(s_2, a) = Q_k^\pi(s_2, a) + \alpha[r' + \gamma Q_k^\pi(s_3, a') - Q_k^\pi(s_2, a)] = 15 + [10 + 10 - 15] = 20$$



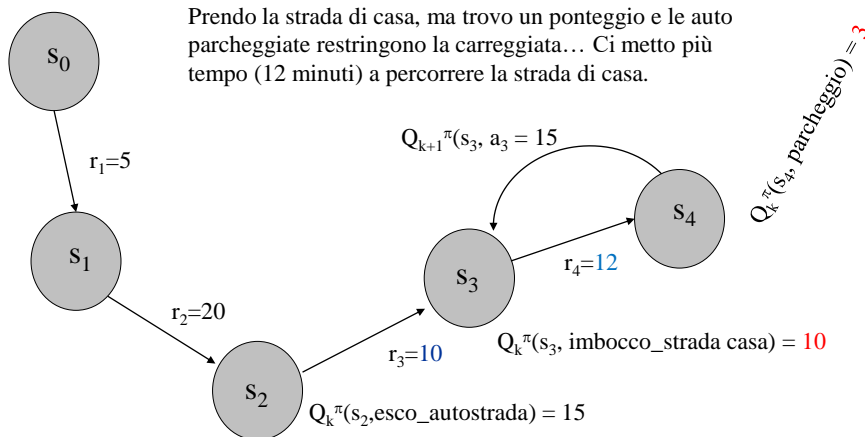
## Learning $Q^\pi(s, a)$ - IV



$s_3 = \text{imbocco\_strada\_casa}; Q_k^\pi(s_3, \text{imbocco\_strada\_casa}) = 10 \text{ min};$

$Q_{k+1}^\pi(s_3, \text{imbocco\_strada\_casa}) = 15 \text{ min};$

$s_0 = \text{ufficio}; Q_{k+1}^\pi(s_0, \text{vado\_parcheggio}) = 43 \text{ minuti}$



Aggiorno il tempo totale dallo stato  $s_2$ :

$$A. Q_{k+1}^\pi(s_2, a) = Q_k^\pi(s_2, a) + \alpha[r' + \gamma Q_k^\pi(s_3, a') - Q_k^\pi(s_2, a)] = 10 + [12 + 3 - 10] = 15 \quad \text{i.it}$$



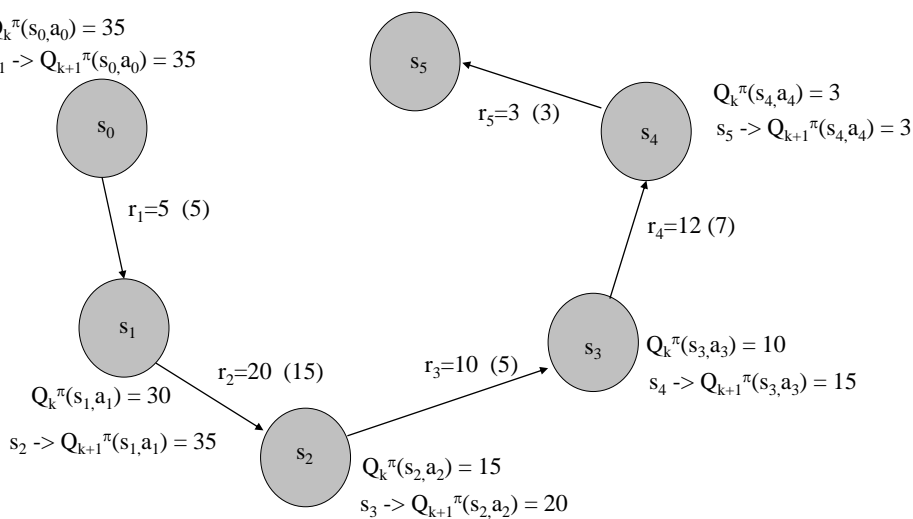
## Learning $Q^\pi(s, a)$



$s_0 = \text{ufficio}; s_5 = \text{casa}.$

$Q_k^\pi(s_0, a_0) = 35$

$s_1 \rightarrow Q_{k+1}^\pi(s_0, a_0) = 35$



Come i diversi reward istantanei modificano  $Q^\pi(s, a)$ ?



## Ruolo di $\alpha$



$$Q_{k+1}(s_1, a_1) = Q_k(s_1, a_1) + \alpha (r_1 + \gamma Q(s_2, a_2) - Q(s_2, a_2)) = 30 + \alpha (20 + 15 - 30) = 30 + \alpha * 5$$

Stima iniziale del tempo di percorrenza dal parcheggio: 30m

Tempo per percorrere l'autostrada: 20m

Stima del tempo di percorrenza dall'uscita del parcheggio: 35min (per  $\alpha = 1$ )

$\alpha < 1$ .

If  $\alpha \ll 1$  aggiorno molto lentamente la value function.

If  $\alpha = 1/k(s, a)$  aggiorno la value function in modo da tendere al valore atteso. Devo memorizzare le occorrenze della coppia stato-azione s, a.

If  $\alpha = \text{cost}$ . Aggiorno la value function, pesando maggiormente i risultati collezionati dalle visite dello stato più recenti.



## Sommario



Le equazioni di Bellman

Differenze temporali