



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

Barto and Sutton, Capitoli 3 e 6



A.A. 2025-2026 1/46 <http://borgheze.di.unimi.it/>

1



World RL competition



<https://neurips.cc/Conferences/2024/CompetitionTrack>
<https://neurips.cc/virtual/2025/loc/san-diego/events/Competition>

Started in NIP2006.

Many competitions, some are based on RL

Other competitions are organized

A.A. 2025-2026 2/46 <http://borgheze.di.unimi.it/>

2



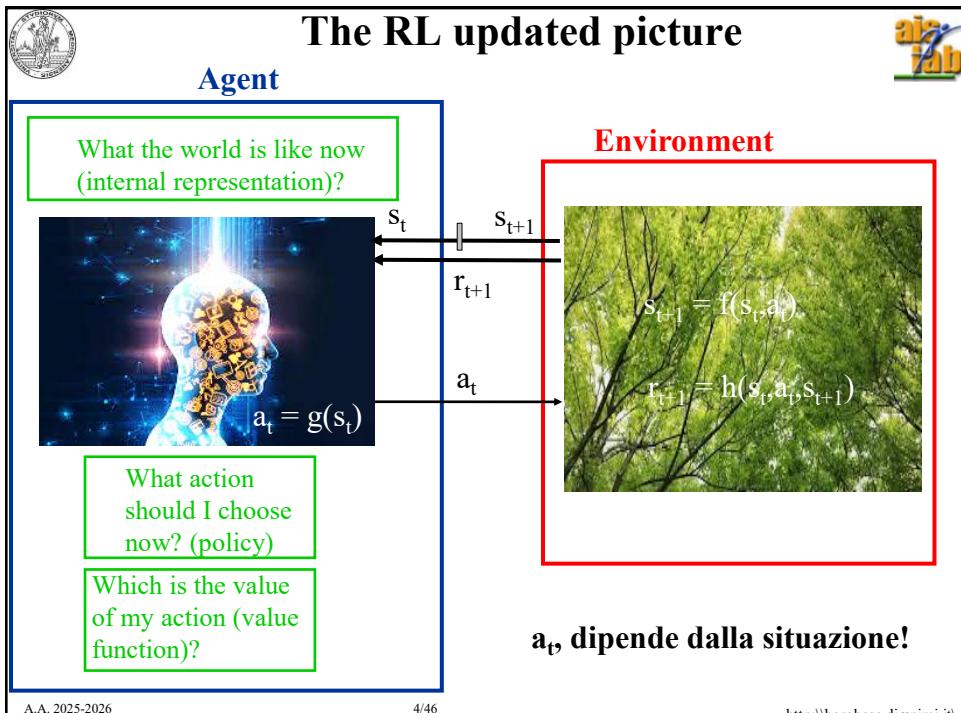
Sommario



Le equazioni di Bellman

Differenze temporali

SARSA





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.

Fino a quando la policy non è stabile.

Policy iteration.



Apprendimento



Iterative policy evaluation

$$Q_{k+1}^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s,s',a} + \gamma \pi(s', a') \sum_{a'} Q_k^\pi(s'_{t+1}, a'_{t+1}) \right\}$$

Converge al limite a $Q^\pi(s,a)$.

Policy improvement

$$Q_{k+1}^\pi(s_t, a_t) = \max_{a'} \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s,s',a} + \gamma \pi(s', a') \sum_{a'} Q_k^\pi(s', a') \right\}$$

Si arriva a:

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



Un ciclo di interazione



Agent

What the world is like now
(internal representation)?



What action
should I choose
now? (policy)

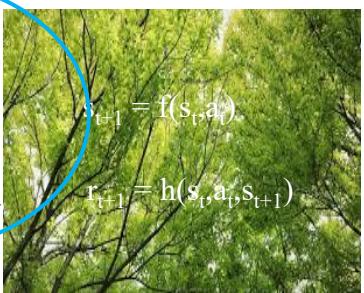
Which is the value
of my action (value
function)?

s_t

r_{t+1}

a_t

Environment



L'agente ragiona su quantità note
a 1 passo

A.A. 2025-2026

7/46

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

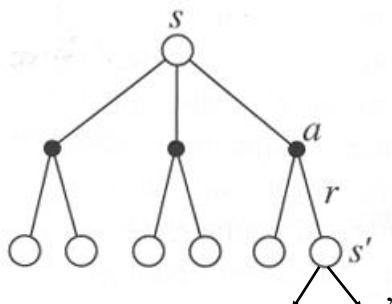
7



Tecnica full-back

Back-up

$\pi(s,a)$ fissata



$t+1$

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 $Q(s, a)$.

π fissata

A.A. 2025-2026

8/46

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

8



Meccanismo di apprendimento nel RL



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

- 1) Implemento una policy ($\pi(s,a)$)
- 2) Stimo la Value function ($Q^\pi(s,a)$) per tutte le coppie stato-azione
- 3) Miglioro la policy, $\pi(s,a)$.

Framework analizzati per calcolare $Q^\pi(s,a)$:

1. Stocastico completo. **L'agente conosce la statistica della dinamica dell'ambiente e dei reward**, scrive **le equazioni lineari** che legano le value function in stati diversi e **calcola i valori di $Q^\pi(.)$** .
2. Stocastico con aggiornamento. **L'agente conosce la statistica della dinamica dell'ambiente e dei reward**. Procede da uno stato iniziale a quello finale. **Da ogni stato esplora in parallelo tutti i possibili stati prossimi** e aggiorna i valori delle $Q^\pi(.)$. Iterative policy evaluation – Policy iteration.
3. Stocastico con interazione singola e con la scelta di una singola azione. **L'agente NON conosce la statistica della dinamica dell'ambiente e dei reward**. Procede da uno stato iniziale a quello finale. Da ogni stato esplora una sola azione e **un SOLO solo stato prossimo**. Aggiorna i valori di $Q^\pi(.)$.



Sommario



Le equazioni di Bellman

Differenze temporali

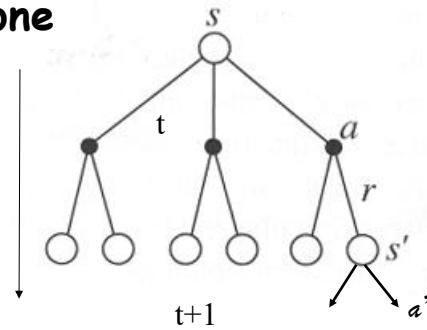
SARSA



1 passo di interazione

Back-up

$\pi(s,a)$ fissata



Analizziamo la transizione da $\{s_t, a_t\} \rightarrow \{s'_{t+1}, a'_{t+1}\}$. Da $\{s', a'\}$ mi guardo indietro e aggiorno $Q(s,a)$.

- Nell'approccio **iterative policy evaluation**, considero TUTTI gli s'_{t+1} (tutti i reward a 1 passo con la probabilità di sceglierli, e i reward a lungo termine da s'_{t+1} : $Q(s'_{t+1}, a'_{t+1}) \forall s'$). *Full backup* se esaminiamo tutti gli s' e a' (cf. Dynamic Programming).

Come estendiamo al caso di un agente che esplora l'ambiente non noto? Non conosciamo tutti gli s'_{t+1} ma 1 SOLO s' , s'_{t+1} , e misuro 1 SOLO reward a un passo; considero il reward TOTALE, a lungo termine, SOLO di 1 stato prossimo.

A.A. 2025-2026

11/46

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

11

$Q(s,a)$ - Osservazioni



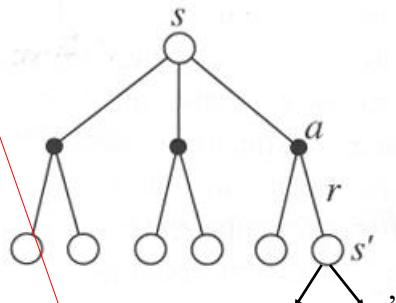
$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s'|a} \left\{ R_{s,s',a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Policy nota

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

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

Cosa cambia?



Non è noto il funzionamento dell'ambiente (interazione)

A.A. 2025-2026

12/46

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

12



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 **se non conosciamo l'ambiente?**



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:

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

$$StepSize = \alpha = 1/(N+1)$$

$$Rewards\ weight\ w = I$$

$$a = cost$$

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

Qual è la differenza introdotta dall'approccio che prevede comportamenti (catene di azioni) e la valutazione di un reward TOTALE, a lungo termine?

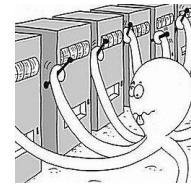


Un possibile aggiornamento di $Q(s, a)$



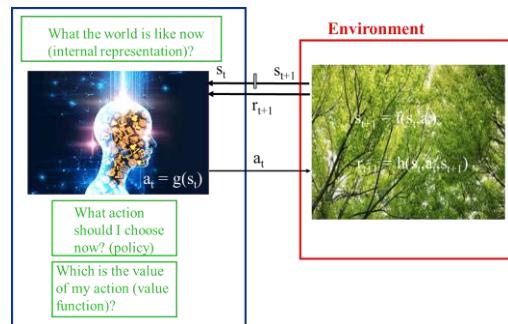
$$Q_{k+1}(a) = Q_k(a) - \frac{Q_k(a)}{N_{k+1}(a)} + \frac{r_{k+1}(a)}{N_{k+1}(a)} = Q_k(a) + \alpha[r_{k+1}(a) - Q_k(a)] =$$

$$Q_{k+1}(a) = Q_k(a) + \alpha \Delta Q_k(a) \quad \alpha = \frac{1}{N_{k+1}(a)}$$



Come passo ai comportamenti?

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



Come calcolo $\Delta Q_k(s, a)$?

A.A. 2025-2026

15/46

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

15



Calcolo di $\Delta Q_k(s, a)$



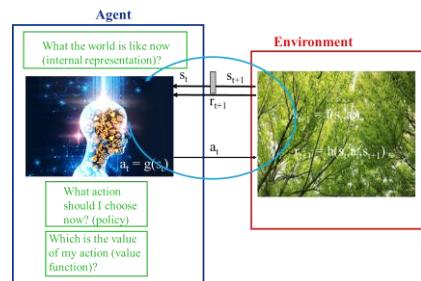
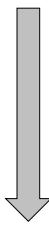
$$Q_{k+1}(a) = Q_k(a) + \alpha[r_{k+1}(a) - Q_k(a)] =$$

$$Q_{k+1}(a) = Q_k(a) + \alpha \Delta Q_k(a)$$

Al tempo t abbiamo a disposizione:

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

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



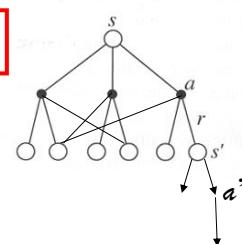
Quale semantica hanno $Q(s, a)$ e $r(s, a, s')$ nel caso dei comportamenti?

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

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

Reward a 1 passo

Reward a lungo termine da s'



A.A. 2025-2026

16/46

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

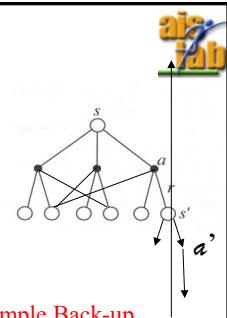
16



TD(0) update

Ad ogni istante di tempo di ogni trial aggiorno 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 $\{s, a\}$: $Q(s_t, a_t)$ congruente con: $Q(s_t, a_t)$ 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 $Q(s, a)$.

Percorro un solo ramo dell'albero, alla volta.

Per α che diminuisce con l'apprendimento, per $k \rightarrow \infty$, idealmente come $1/(k+1)$,
 $Q_k^{\pi}(s, a)$ converge al valore vero di $Q^{\pi}(s, a)$

$\pi(s, a)$ fissata

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



Confronto con il setting associativo



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

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:

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

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



Setting α 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*.

α 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 definito (policy 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:

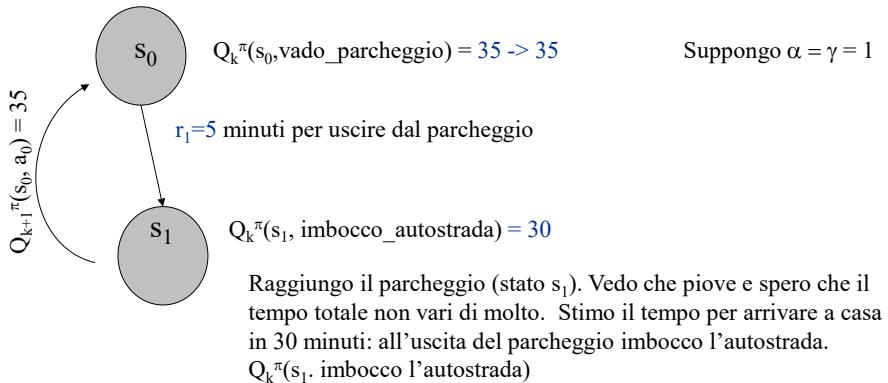
- Dall'ufficio (time to go = 35 minuti) - partito
- Dall'ufficio all'uscita del parcheggio: 5 minuti (time to go = 30 minuti)
- Dal parcheggio all'uscita dell'autostrada: 15 minuti (time to go = 15 minuti)
- Dall'uscita dell'autostrada alla strada di casa: 5 minuti (time to go = 10 minuti)
- Dalla strada di casa al parcheggio di casa: 7 minuti (time to go = 3 minuti)
- Dal parcheggio a casa: 3 minuti (time to go = 0 minuti) - arrivato
- In totale 35 minuti.



Learning $Q^\pi(s_0, a_0)$ - I



s_0 = ufficio; $Q_k^\pi(s_0, \text{vado_parcheggio}) = 35$ 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).



Ricalcolo il tempo totale, ovverosia il tempo dallo stato s_0 che non varia:

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

A.A. 2025-2026

21/46

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

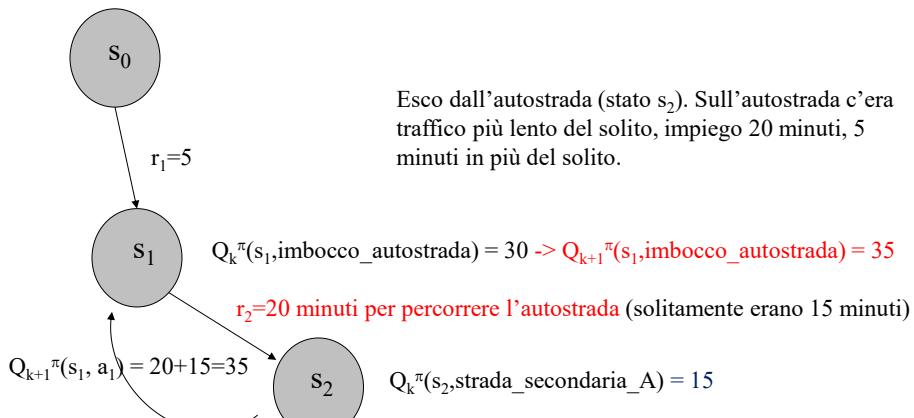
21



Learning $Q^\pi(s_1, a_1)$ - II



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



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

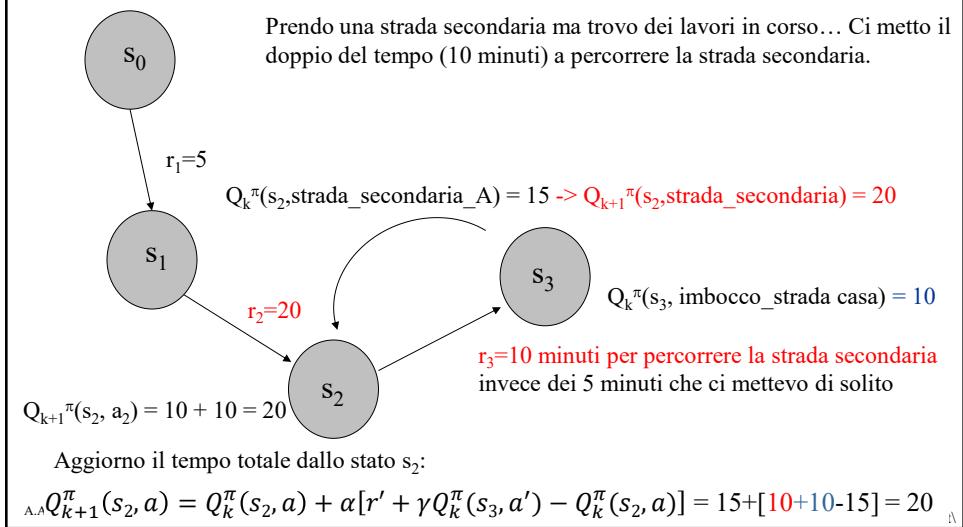
22



Learning $Q^\pi(s_2, a_2)$ - III



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



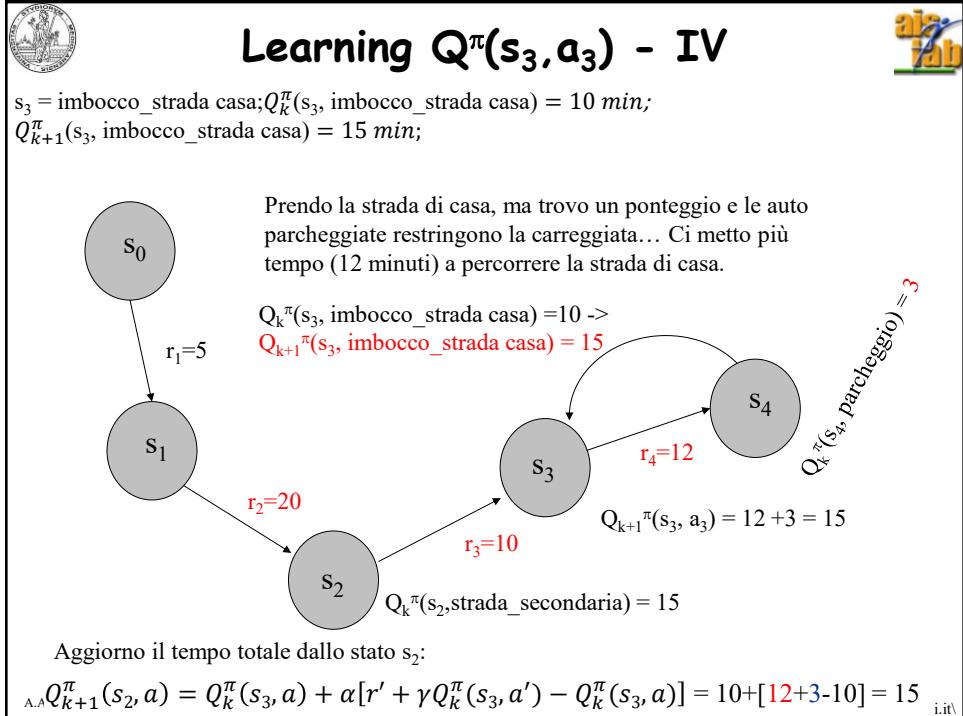
23



Learning $Q^\pi(s_3, a_3)$ - IV



$s_3 = \text{imbocco_strada casa}; Q_k^\pi(s_3, \text{imbocco_strada casa}) = 10 \text{ min}; Q_{k+1}^\pi(s_3, \text{imbocco_strada casa}) = 15 \text{ min};$



24



Esempio - dopo il primo trial



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

La durata dei diversi segmenti può variare da giorno a giorno e quindi la stima della durata totale è stata correttamente conseguentemente all'esplorazione.

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

- Dall'ufficio all'uscita del parcheggio: 5 minuti (time to go = 35 minuti) - incongruenza
- Dal parcheggio all'uscita dell'autostrada: 15 minuti (time to go = 35 minuti)
- Dall'uscita dell'autostrada alla strada di casa: 5 minuti (time to go = 20 minuti)
- Dalla strada di casa a casa: 7 minuti (time to go = 15 minuti)
- Dal parcheggio a casa: 3 minuti (time to go = 3 minuti)

Si sono create diverse incongruenze (ad esempio il time to go è di 35 minuti dall'ufficio come dal parcheggio!), che verranno corrette via via che si ripeteranno le stesse situazioni.

Attualmente la stima aggiornata di $Q(\cdot)$ è per lo stato prima di quello finale ed è di 15 minuti. La stima di $Q(\cdot)$ per gli stati precedenti, viene via via aggiornata nei trial successivi.

In totale ci metto 50 minuti. Come i diversi reward istantanei modificano $Q^\pi(s,a)$ per tutti gli stati?



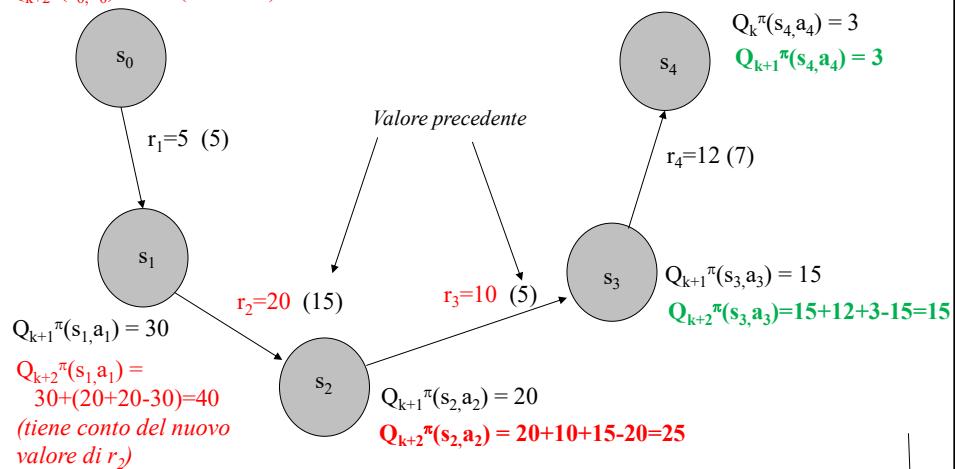
Learning $Q^\pi(s,a)$ - Trial 2

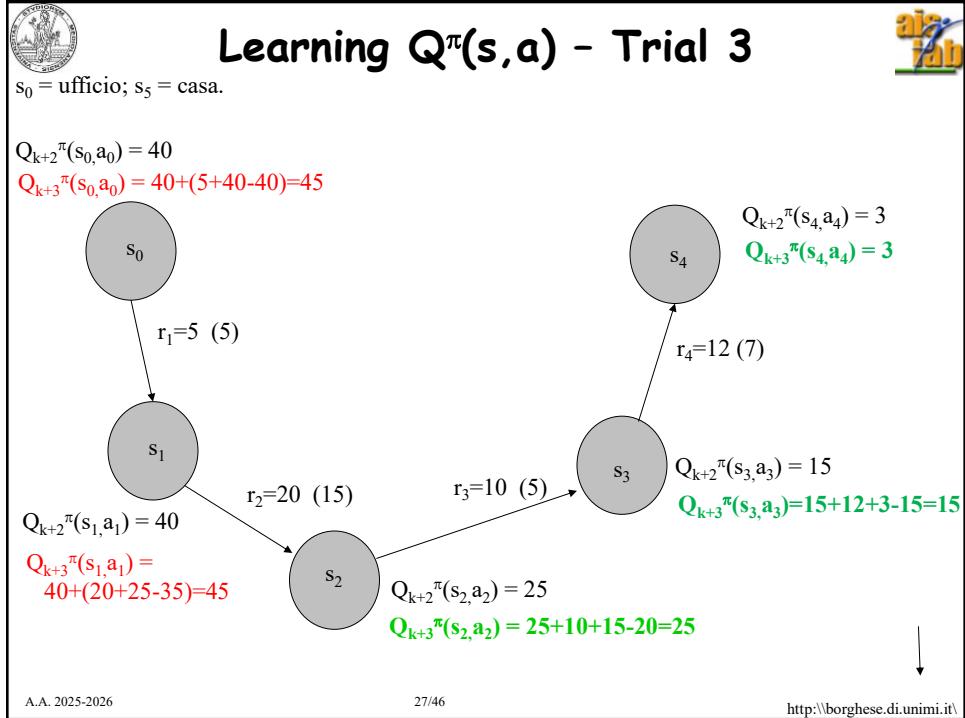


s_0 = ufficio; s_5 = casa.

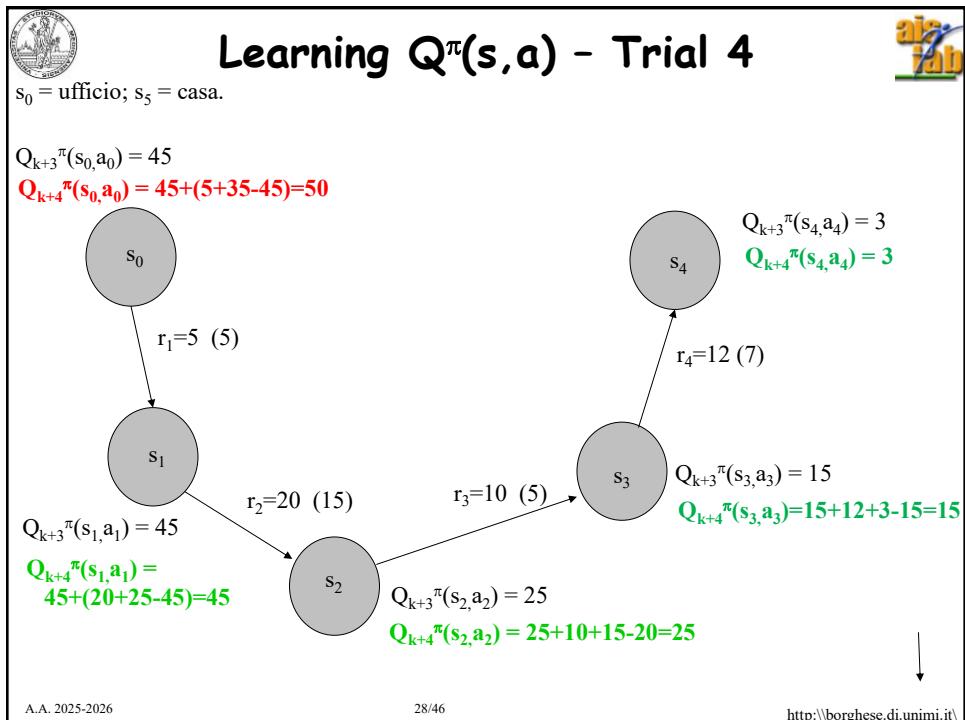
$$Q_{k+1}^\pi(s_0, a_0) = 35$$

$$Q_{k+2}^\pi(s_0, a_0) = 35 + (5 + 35 - 35) = 40 \text{ - corretto}$$





27



28



Osservazioni



L'apprendimento procede dalla fine verso l'inizio!

Non posso sapere all'inizio come si concluderà l'episodio, a ogni passo devo cercare di ottenere una stima migliore del reward che otterrò in futuro (per potere poi scegliere)

La stima di $Q^\pi(s_0, a_0)$ migliora con i trial:

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

$$Q_{k+1}^\pi(s_0, a_0) = 40$$

$$Q_{k+2}^\pi(s_0, a_0) = 45$$

$$Q_{k+3}^\pi(s_0, a_0) = 50$$

E migliora con il miglioramento della stima del valore dello stato successivo: $Q^\pi(s_1, a_1)$



Ruolo di α



$(\alpha < 1)$

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

Stima iniziale del tempo di percorrenza dal parcheggio: 35 m

Tempo per percorrere l'autostrada: 20m

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

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.

La convergenza è garantita per α che decresce gradualmente verso zero.



Sommario



Le equazioni di Bellman

Differenze temporali

SARSA



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):

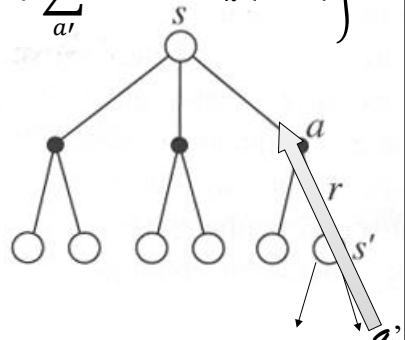
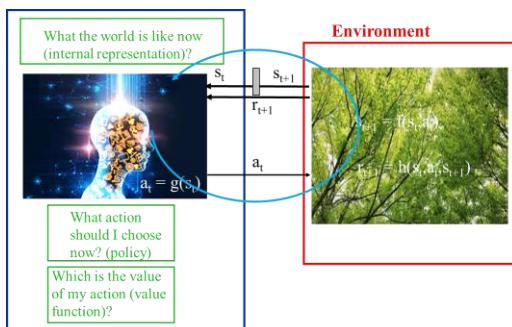
- 1) Implemento una policy ($\pi(s,a)$)
- 2) Calcolo la Value function ($Q^\pi(s,a)$)
- 3) Aggiorno la policy.



Un ciclo di interazione

$$Q_{k+1}^{\pi}(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \sum_{a'} \pi(s', a') Q_k^{\pi}(s', a') \right\}$$

Calcolo ricorsivo di Q(.)



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



SARSA

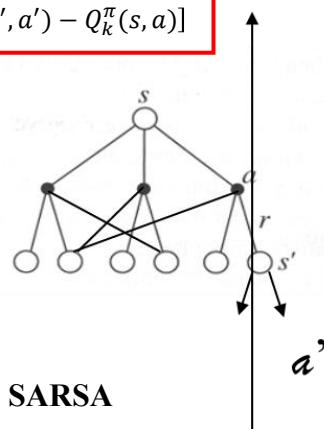
Non richiede conoscenze a priori dell'ambiente.

L'agente stima a partire da nulla (bootstrap).

Si dimostra che il metodo **converge asintoticamente**, stima $Q^{\pi}(s, a)$ quando α decresce dolcemente a 0 all'aumentare del numero di trial

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

Sample backup, single state,
 s_p , single action, a_p , single
future state $s' = s_{t+1}$



State-Action-Reward-State-Action => SARSA



Background su SARSA



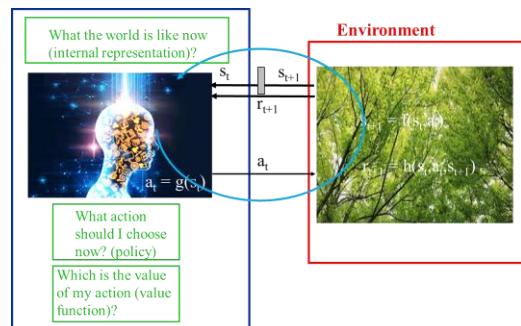
Al tempo t abbiamo a disposizione:

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

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

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

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



A.A. 2025-2026

35/46

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

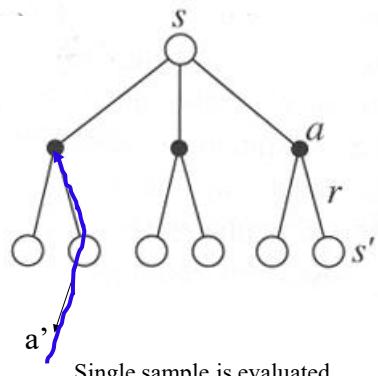
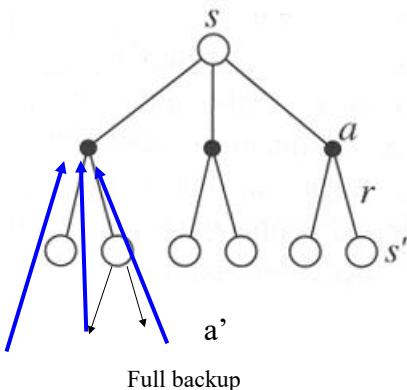
35



Sample backup



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



State = s, Action = a, Reward (a un passo) = r, State = s', action = a'

A.A. 2025-2026

36/46

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

36



SARSA Algorithm (progetto)



```

Q(s,a) = rand();           // ∀s, ∀a, eventualmente Q(s,a) = 0 - inizializzazione
Policy definita;          // Policy specificata, eventualmente stocastica
Repeat                      // for each episode
{
  s = s0;
  Repeat                    // for each step of the actual episode
  {
    a = Policy(s);          // policy deterministica o stocastica
    snext = NextState(s,a); // Funzione non nota all'agente
    reward = Reward(s,snext,a);
    anext = Policy(snext);
    Q(s,a) = Q(s,a) + α [reward + γ Q(snext, anext) − Q(s,a)];
    s = snext;
  }                         // until last state
}                         // until the end of learning (convergence of Q(s,a) to true Q(s,a) ∀s, ∀a, for policy π(s,a) )

```

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

Come integrare i due passi?

A.A. 2025-2026

37/46

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

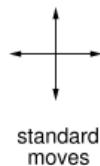
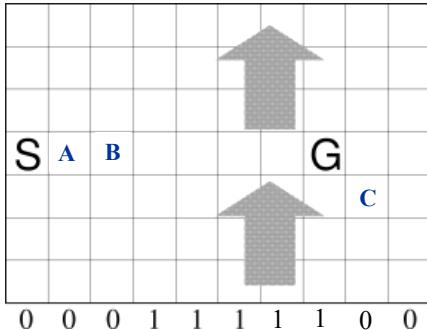
37



Esempio



From Start, S, to Goal, G.



Stati = {caselle della griglia}
 Stato iniziale = S
 Terminal state = G
 Azioni = {su, destra, giù, sinistra}
 Reward = -1 tranne che quando s' = TS (reward = 0)

Upwards wind: somma uno spostamento verso l'alto allo spostamento dell'agente nelle prime sei righe del labirinto.

A.A. 2025-2026

38/46

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

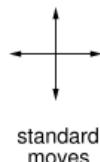
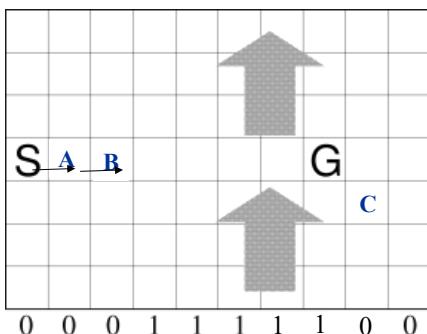
38



Esempio



From Start, S, to Goal, G.



$$\alpha = 0.5 \\ \gamma = 1$$

Upwards wind

$Q(s,a)$ iniziale = 0.

$r = 0$ se $s' = G$; altrimenti $r = -1$.

$\pi(s,a)$ data. Azioni possibili: N, S, E, W, NE, NO, SE, SO

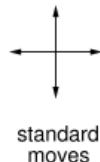
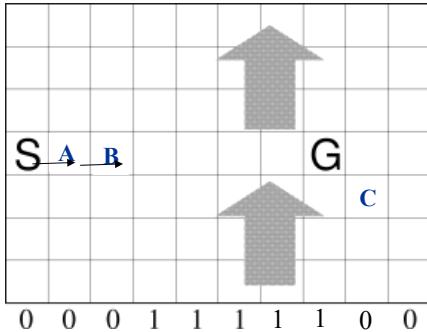
$$\text{Vogliamo valutare } \pi(s,a). \quad Q_{k+1}^{\pi}(s, a) = Q_k^{\pi}(s, a) + \alpha[r' + \gamma Q_k^{\pi}(s', a') - Q_k^{\pi}(s, a)]$$



Esempio - risultato



$$\alpha = 0.5 \\ \gamma = 1$$



Correzione di Q ad un passo:

$$Q_{k+1}^{\pi}(S, east) = Q_k^{\pi}(S, east) + \alpha[r' + \gamma Q_k^{\pi}(A, east) - Q_k^{\pi}(S, east)] = 0 + 0.5 [-1 + 1 \times 0 - 0] = -0.5$$

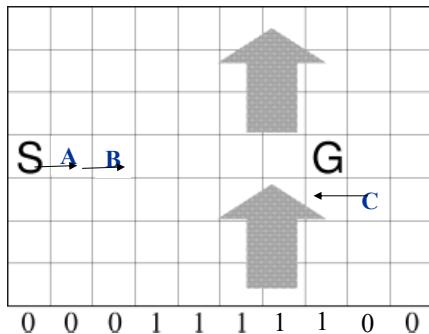
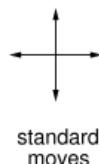
$$Q_{k+1}^{\pi}(s, a) = Q_k^{\pi}(s, a) + \alpha[r' + \gamma Q_k^{\pi}(s', a') - Q_k^{\pi}(s, a)]$$



Esempio - risultato



$$\alpha = 0.5 \\ \gamma = 1$$



$$Q_{k+1}^{\pi}(A, east) = -0.5$$

Correzione di Q ad un passo:

$$Q_{k+2}^{\pi}(S, east) = Q_{k+1}^{\pi}(S, east) + \alpha[r' + \gamma Q_{k+1}^{\pi}(A, east) - Q_{k+1}^{\pi}(S, east)] = \\ -0.5 + 0.5 [-1+1x(-0.5) - 0.5] = -1$$

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

A.A. 2025-2026

41/46

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

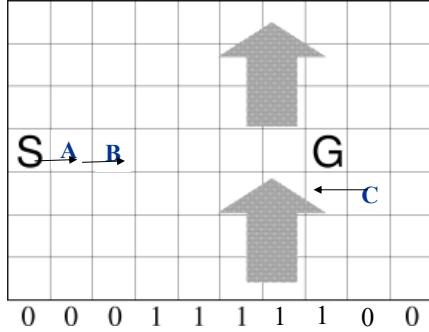
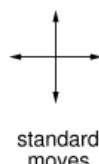
41



Esempio - risultato



$$\alpha = 0.5 \\ \gamma = 1$$



Correzione di Q ad un passo:

$$Q_{k+1}^{\pi}(A, east) = Q_k^{\pi}(A, east) + \alpha[r' + \gamma Q_k^{\pi}(B, east) - Q_k^{\pi}(A, east)] = 0 + 0.5 [-1+1x0-0] = -0.5$$

$$Q_{k+1}^{\pi}(C, west) = Q_k^{\pi}(C, west) + \alpha[r' + \gamma Q_k^{\pi}(G, .) - Q_k^{\pi}(C, west)] = 0 + 0.5 [0+1x0-0] = 0 \\ (\text{NB c'è il vento verso l'alto di 1})$$

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

A.A. 2025-2026

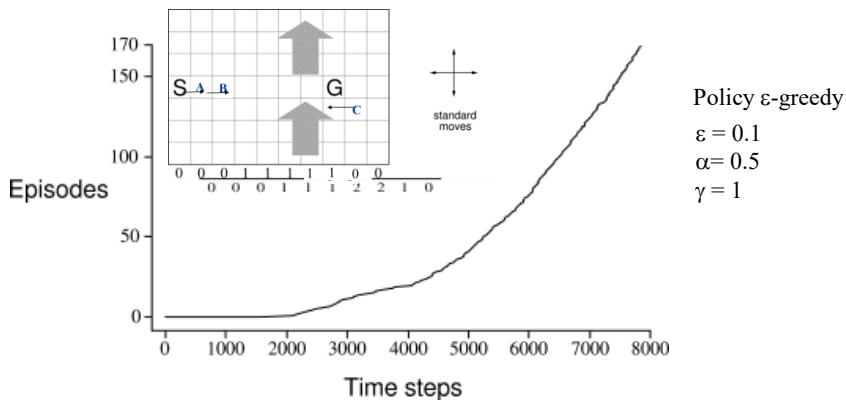
42/46

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

42



Esempio - risultato

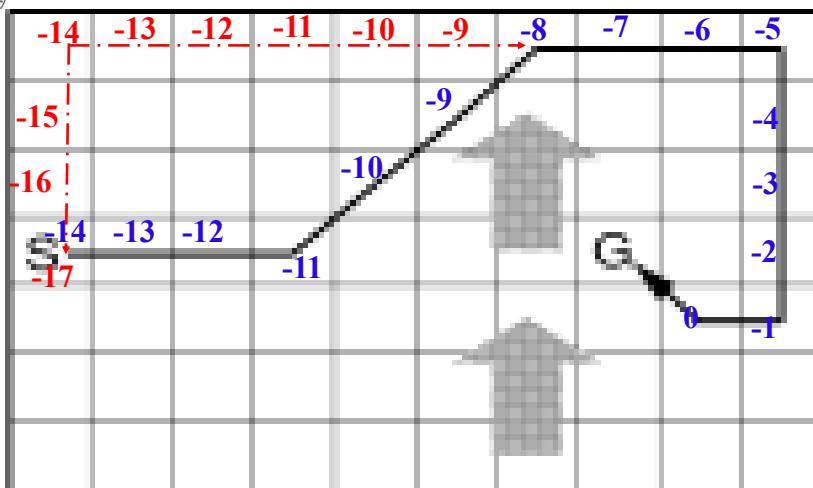


Aumentano gli episodi nello stesso intervallo di tempo, via via che trial vengono eseguiti.
All'inizio un trial richiede molto tempo per essere eseguito.

Non è il percorso ottimo (17 passi contro 15 passi)
E' il percorso cristallizzato.



Esempio - valore di $Q^\pi(\cdot)$

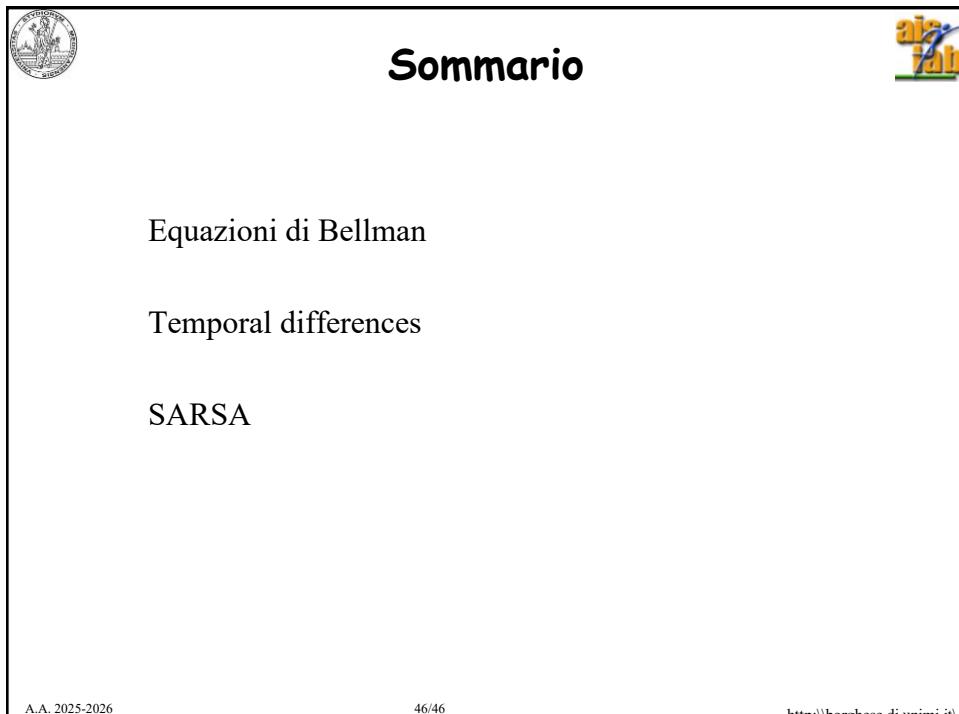


We have estimated $Q^\pi(s, a)$. In most states we have better option.

E.g. in S if we have better $Q^\pi(S, a_{new})$ than moving East. This allows exploring different paths. We can change after having estimated $Q^\pi(\cdot)$.



45



46