

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

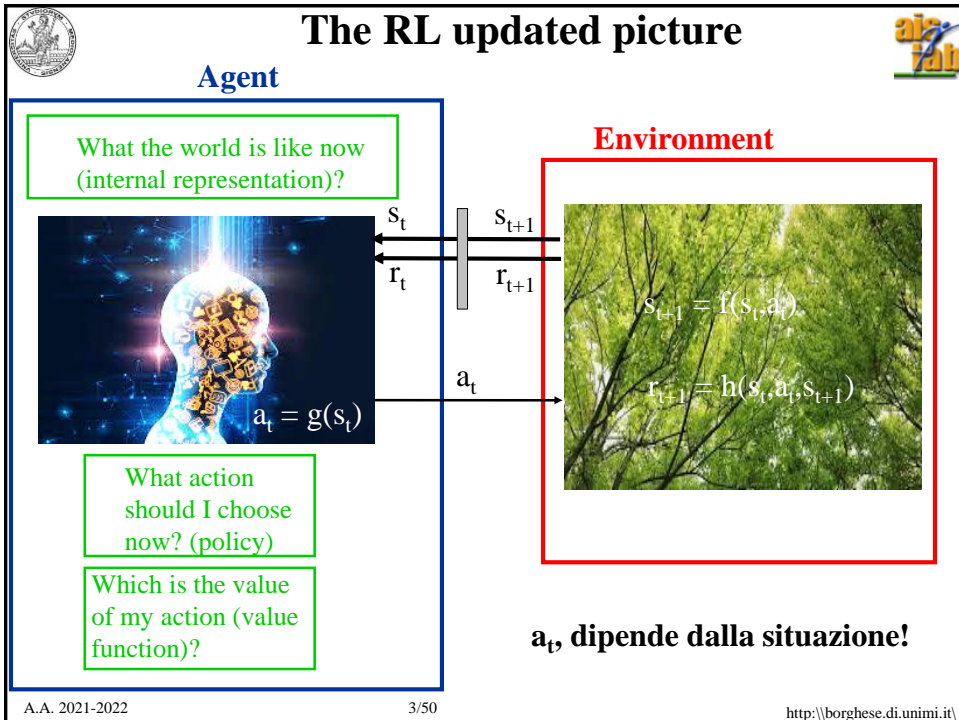


Sommario



Le equazioni di Bellman

Differenze temporali



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. 2021-2022

4/50

<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. 2021-2022

5/50

<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

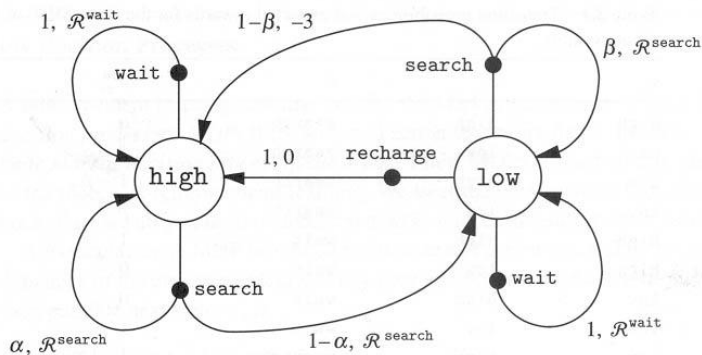
$Q(\text{high, search}) = ?$

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

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

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

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





A.A. 2021-2022

6/50

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

Analisi ad un passo dal tempo t



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

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

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

A.A. 2021-2022 7/50 http://borghese.di.unimi.it/

Analisi ad un passo dal tempo t

Policy deterministica
 a(high) = wait
 a(low) = 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. 2021-2022 8/50 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. 2021-2022 9/50 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!!

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

2 cammini possibili!!

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

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

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

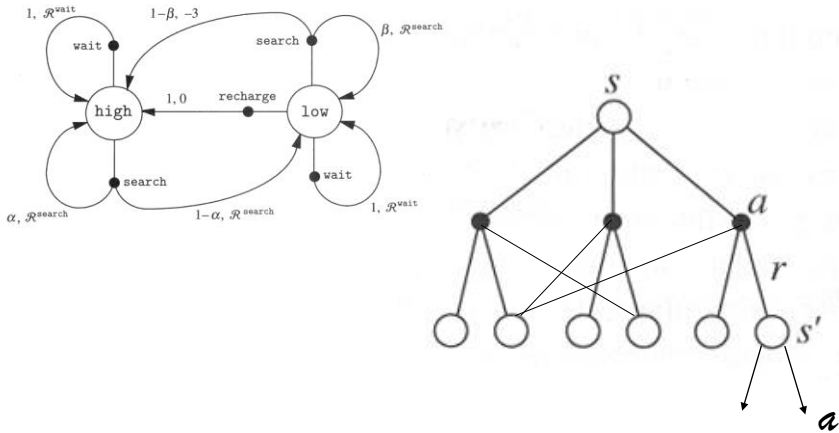
Contiene la probabilità di ricevere un reward $\gamma Q(s', a)$, condizionata a $s_{t+1} = s'$ http://borghese.di.unimi.it/



Valutazione policy stocastica



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



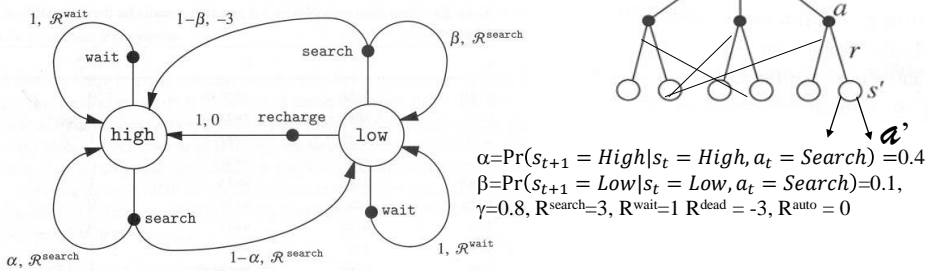
A.A. 2021-2022

11/50

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



Policy stocastica



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

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

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

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

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

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

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

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

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

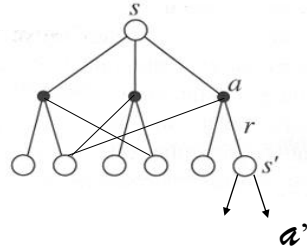
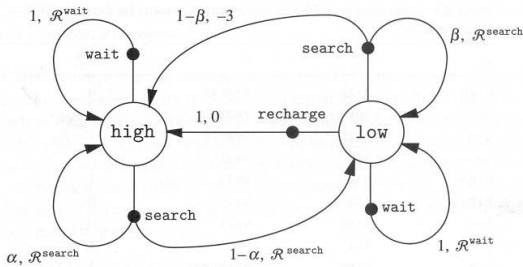
A.A. 2021-2022

12/50

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



Policy stocastica



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

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

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

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

A.A. 2021-2022

5 equazioni in 5 incognite

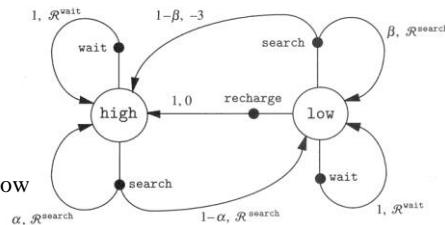
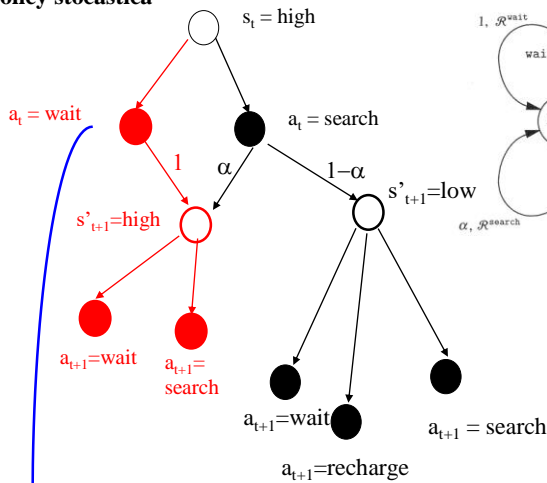
<http://borghese.di.unimi.it/>



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. 2021-2022

14/50

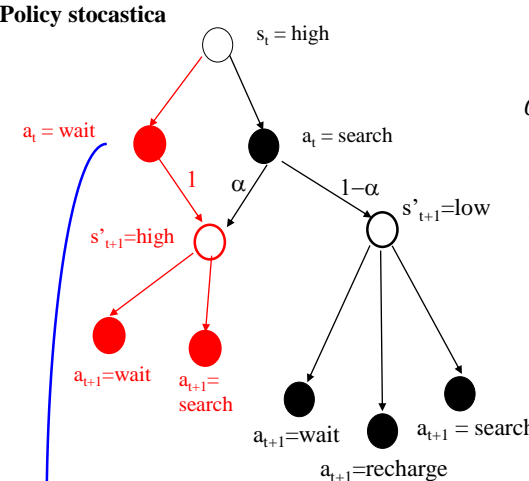
<http://borghese.di.unimi.it/>



Analisi ad un passo dal tempo t



Policy stocastica

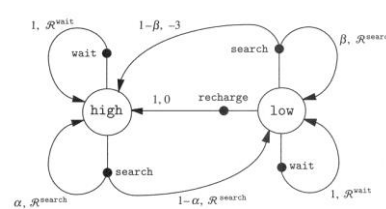


$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$


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

2 cammini possibili!!


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



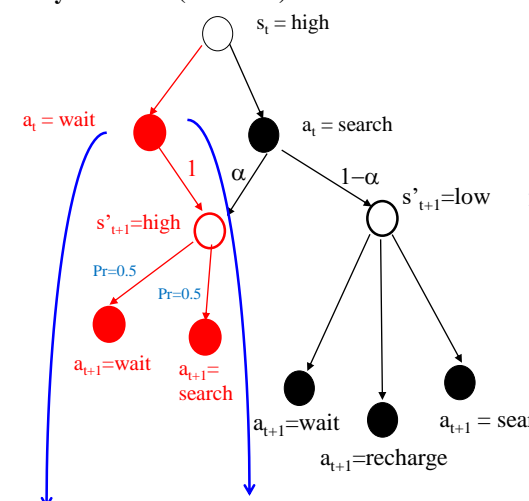
A.A. 2021-2022 15/50 http://borghese.di.unimi.it/



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

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

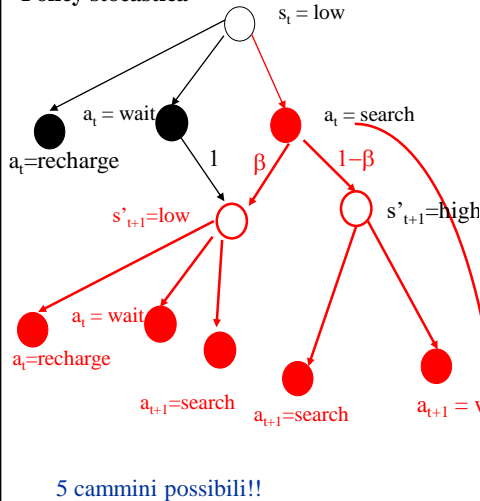
A.A. 2021-2022 16/50 http://borghese.di.unimi.it/



Analisi ad un passo dal tempo t



Policy stocastica



$s_t = \text{low}$

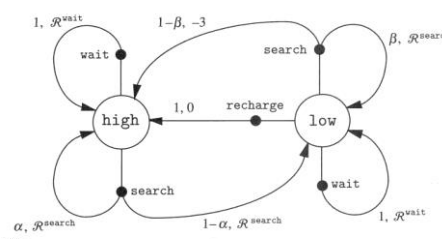
$a_t = \text{wait}$ $a_t = \text{search}$

$a_t = \text{recharge}$

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

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

5 cammini possibili!!




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


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

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

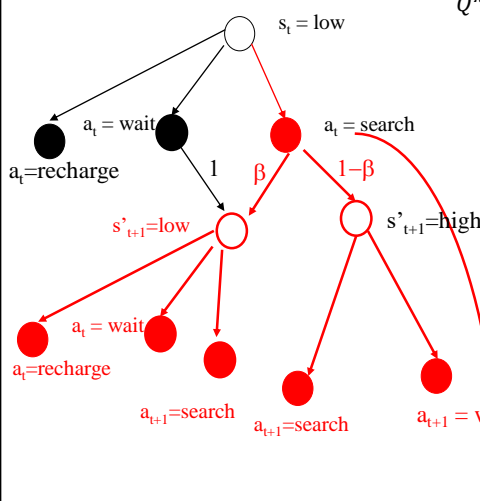
A.A. 2021-2022 17/50 http://borghese.di.unimi.it/



Analisi ad un passo dal tempo t



Policy stocastica (equiprobabile)



$s_t = \text{low}$

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

$a_t = \text{recharge}$

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

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

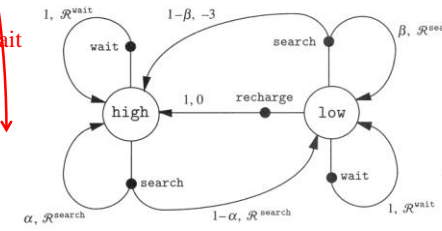
5 cammini possibili!!

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


5 cammini possibili!!

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




A.A. 2021-2022 18/50 http://borghese.di.unimi.it/



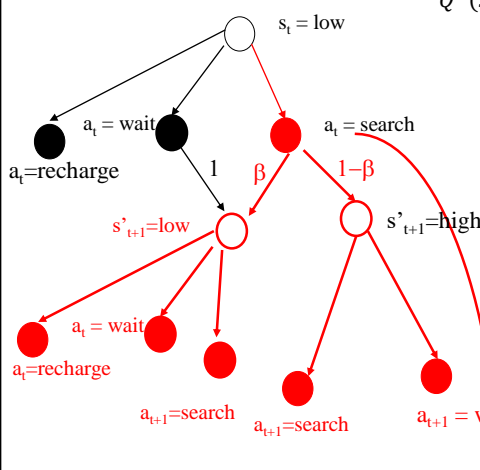
Analisi ad un passo dal tempo t

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. 2021-2022
19/50
<http://borghese.di.unimi.it/>



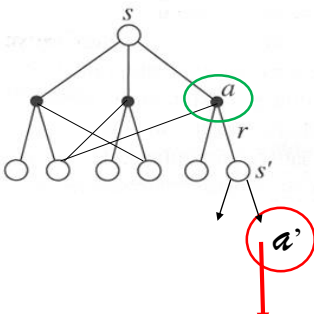
Calcolo ricorsivo della Value function



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

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

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



A.A. 2021-2022
20/50
<http://borghese.di.unimi.it/>



Calcolo ricorsivo della Value function



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

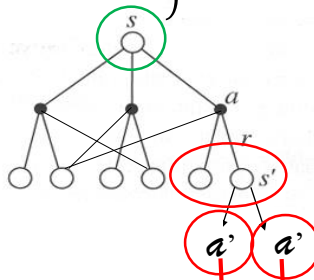
Isolo il reward ad un passo nella serie dei reward.

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

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

Io termine
(a un passo)

Io termine
(passi futuri)



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

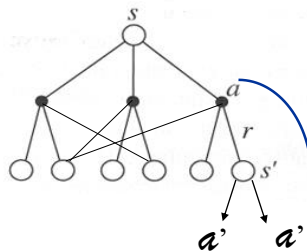


$$P_{s \rightarrow s' | a} \triangleq \Pr(s_{t+1} = s' | s_t = s, a_t = a)$$

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

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' : $R_{s \rightarrow s' | a_j}$

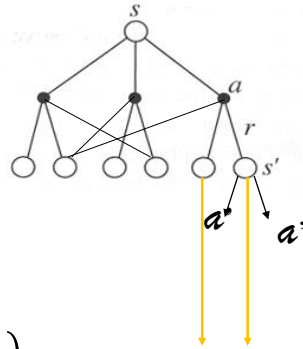


$Q^\pi(s, a)$: secondo termine



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

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



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

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



Putting all together



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

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

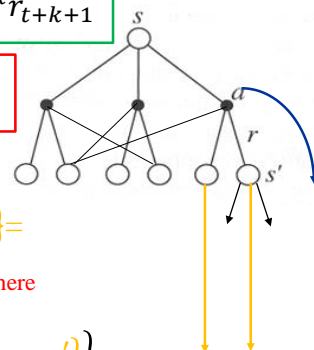
$$\sum_{s'} P_{s \rightarrow s' | a} R_{s, s', a} +$$

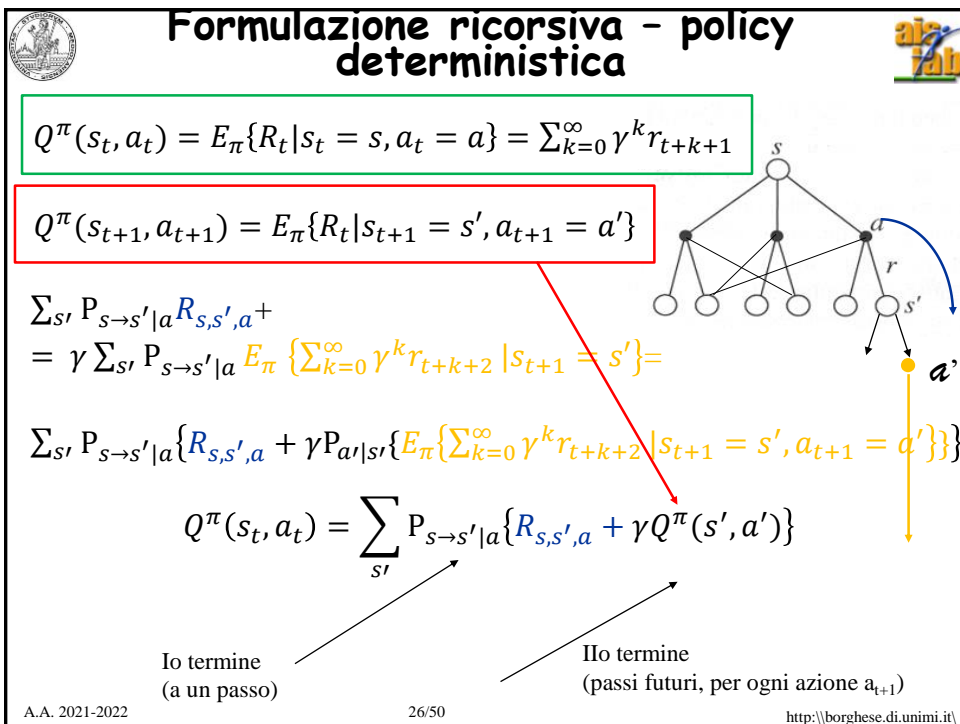
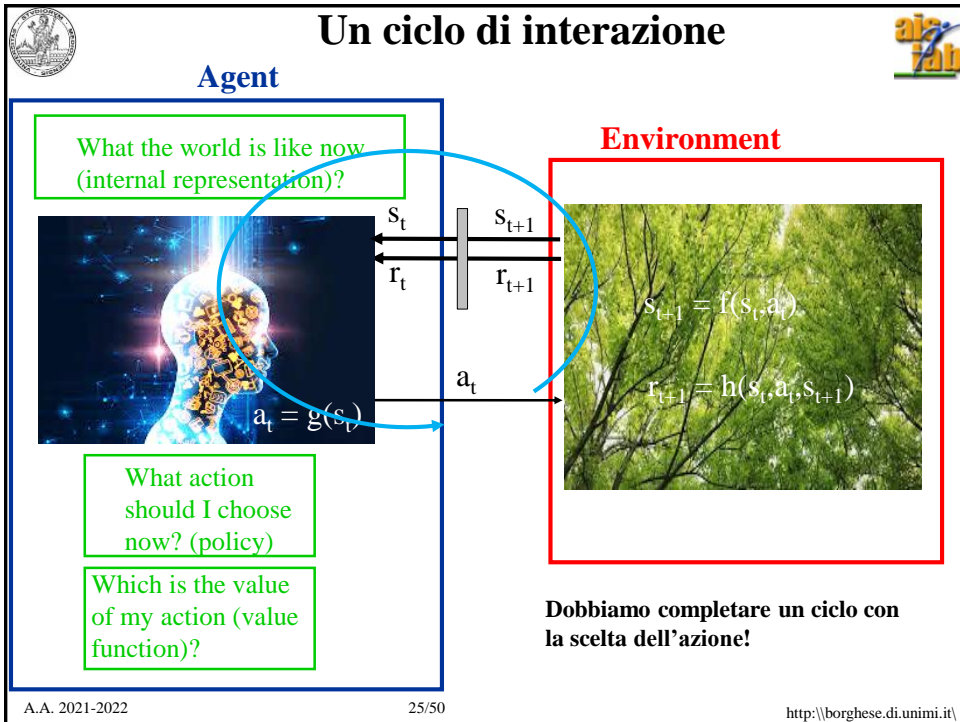
$$= \gamma \sum_{s'} P_{s \rightarrow s' | a} E_\pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_{t+1} = s' \right\} =$$


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

Io termine
(a un passo)


Il termine
(passi futuri)







Formulazione ricorsiva - policy stocastica



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

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

$$\sum_{s'} P_{s \rightarrow s' | a} R_{s, s', a} +$$

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


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

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


Io termine
(a un passo)

Io termine
(passi futuri, per ogni azione a_{t+1})

A.A. 2021-2022 27/50 http://borghese.di.unimi.it/





Equazioni di Bellman



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

Calcolo ricorsivo di Q(.)

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


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(.)

Agent


What the world is like now (internal representation)?



What action should I choose now? (policy)

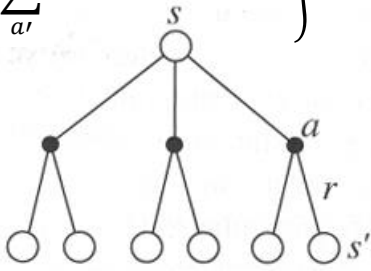
Which is the value of my action (value function)?

Environment





$r = r(s, a)$

$s' = b(s, a, s_{-1})$




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

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

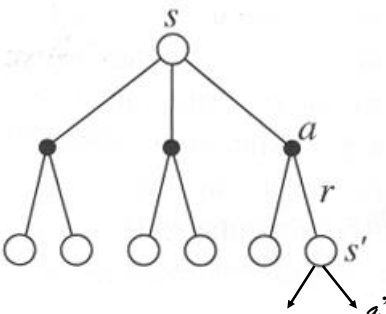



Tecnica full-back

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

<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):

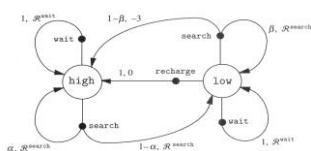
- 1) Implemento una policy ($\pi(s,a)$)
- 2) Aggiorno la Value function ($Q^\pi(s,a)$)
- 3) **Aggiorno la policy.**



$Q(s,a)$ - Osservazioni

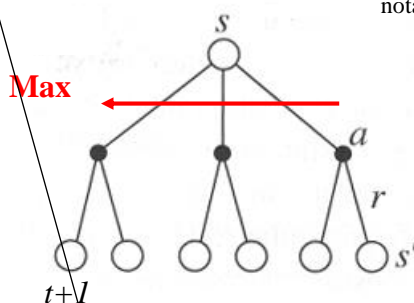


$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s,s',a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$



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



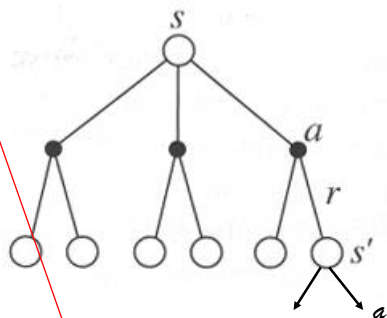
$$Q^\pi(s_t, a_t) = \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \sum_{a'} \pi(s', a') Q^\pi(s', a') \right\}$$

Policy nota

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

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

Cosa cambia?



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



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

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:

NewEstimate = OldEstimate + StepSize [Target - OldEstimate]

*NewEstimate = OldEstimate + StepSize * Error.*

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

$a = cost$

Rewards weight $w = 1$

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)?

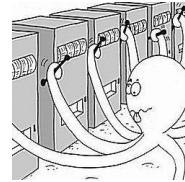


Un possibile aggiornamento di $Q(s, a)$



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

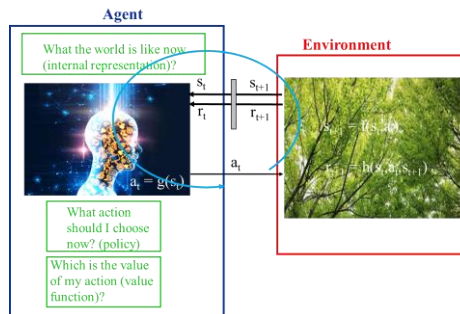
$$Q_{k+1}(a) = Q_k(a) + \alpha \Delta Q_k(a)$$



Come passo ai comportamenti?

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

Come calcolo ΔQ_k ?



Calcolo di ΔQ_k



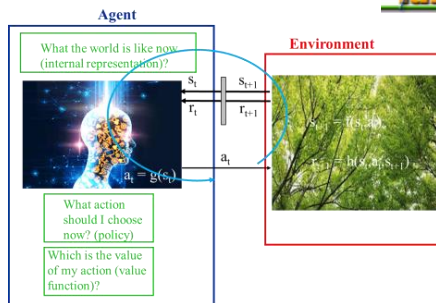
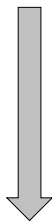
$$Q_{k+1}(a) = Q_k(a) + \alpha[r_{k+1}(a) - Q_k(a)] =$$

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

Al tempo t abbiamo a disposizione:

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

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



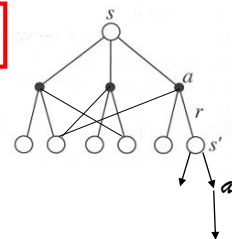
Quale semantica hanno $Q(s, a)$ e $r(s, a, s')$ nel caso dei comportamenti?

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

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

Reward a 1 passo

Reward a lungo termine da s'



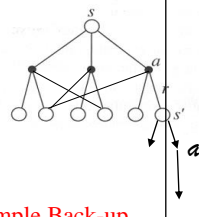


TD(0) update



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 $\{s, a\}$: $Q(s_t, a_t)$ 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 $Q(s, a)$.

Percorro un solo ramo dell'albero, alla volta.

Per α 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 $Q^{\pi}(s, a)$



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:

$$NewEstimate = OldEstimate + StepSize [Target - OldEstimate]$$

$$NewEstimate = OldEstimate + StepSize * Error.$$

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



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 definite (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 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)

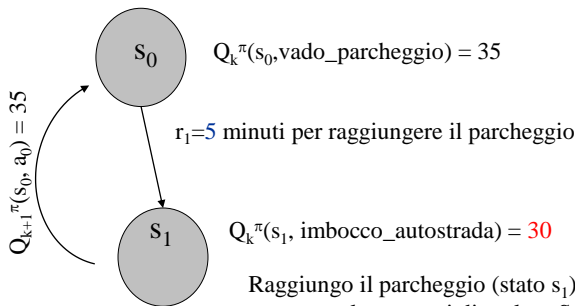
- In totale 35 minuti.



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



s_0 = ufficio; $Q_k^\pi(s_0, \text{vado_parcheggio}) = 35$ minuti; $Q_{k+1}^\pi(s_0, \text{vado_parcheggio}) = 35$ 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).



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

$r_1 = 5$ minuti per raggiungere il parcheggio

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

Raggiungo il parcheggio (stato s_1). Vedo che piove e spero che il tempo totale non vari di molto. Stimo il tempo per arrivare a casa in 30 minuti: all'uscita del parcheggio imbocco l'autostrada $Q_k^\pi(s_1, \text{imbocco_autostrada})$

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 $\alpha = \gamma = 1$

A.A. 2021-2022

43/50

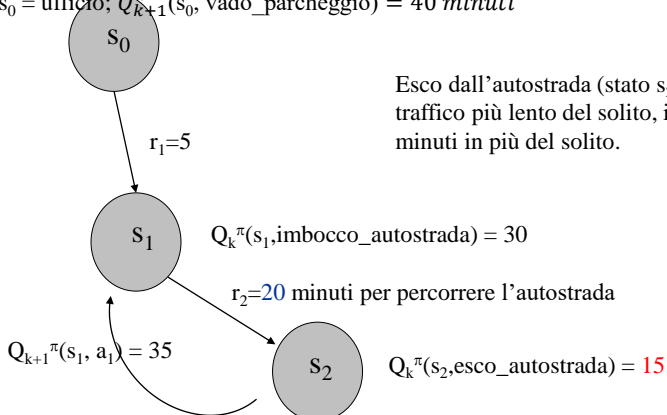
<http://borghese.di.unimi.it/>



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



s_1 = parcheggio; $Q_k^\pi(s_1, \text{imbocco_autostrada}) = 30$ minuti;
 $Q_{k+1}^\pi(s_1, \text{imbocco_autostrada}) = 30$ 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)
 s_0 = ufficio; $Q_{k+1}^\pi(s_0, \text{vado_parcheggio}) = 40$ minuti



Esco dall'autostrada (stato s_2). Sull'autostrada c'era traffico più lento del solito, impiego 20 minuti, 5 minuti in più del solito.

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

$r_2 = 20$ minuti per percorrere l'autostrada

$Q_k^\pi(s_2, \text{esco_autostrada}) = 15$

Aggiorno il tempo totale dallo stato s_1 :

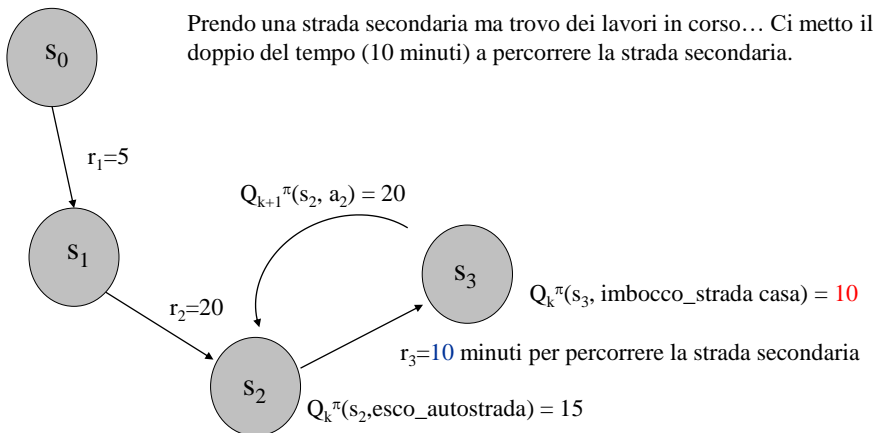
$$Q_{k+1}^\pi(s_1, a) = Q_k^\pi(s_1, a) + \alpha[r' + \gamma Q_k^\pi(s_2, a') - Q_k^\pi(s_1, a)] = 30 + [20 + 15 - 30] = 35$$



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 :

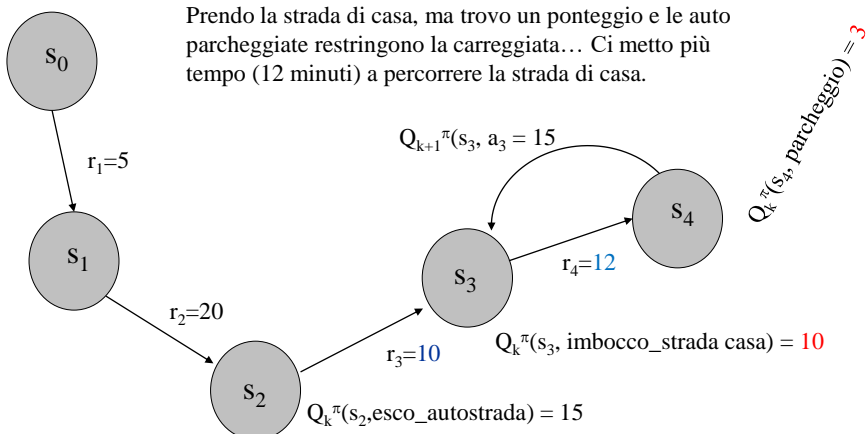
$$A. Q_{k+1}^\pi(s_2, a) = Q_k^\pi(s_2, a) + \alpha[r' + \gamma Q_k^\pi(s_3, a') - Q_k^\pi(s_2, a)] = 15 + [10 + 10 - 15] = 20$$



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 :

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



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



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

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

$Q_k^\pi(s_1, a_1) = 30$
 $s_2 \rightarrow Q_{k+1}^\pi(s_1, a_1) = 35$

$Q_k^\pi(s_2, a_2) = 15$
 $s_3 \rightarrow Q_{k+1}^\pi(s_2, a_2) = 20$

$Q_k^\pi(s_3, a_3) = 10$
 $s_4 \rightarrow Q_{k+1}^\pi(s_3, a_3) = 15$

$Q_k^\pi(s_4, a_4) = 3$
 $s_5 \rightarrow Q_{k+1}^\pi(s_4, a_4) = 3$


Transitions and rewards:
 $s_0 \xrightarrow{r_1=5 (5)} s_1$
 $s_1 \xrightarrow{r_2=20 (15)} s_2$
 $s_2 \xrightarrow{r_3=10 (5)} s_3$
 $s_3 \xrightarrow{r_4=12 (7)} s_4$
 $s_4 \xrightarrow{r_5=3 (3)} s_5$

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


A.A. 2021-2022

47/50

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



Esempio - dopo il primo trial



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

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

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 = 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 3 minuti. La stima di $Q(\cdot)$ per gli stati precedenti, viene via via aggiornata nei trial successivi.

A.A. 2021-2022

48/50

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



Ruolo di α



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

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.

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



Sommario



Le equazioni di Bellman

Differenze temporali