

Sistemi Intelligenti Reinforcement Learning: Policy iteration

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. 2024-2025

1/62

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



Sommario



Le equazioni di Bellman per policy deterministica

Le equazioni Bellman per policy stocastica

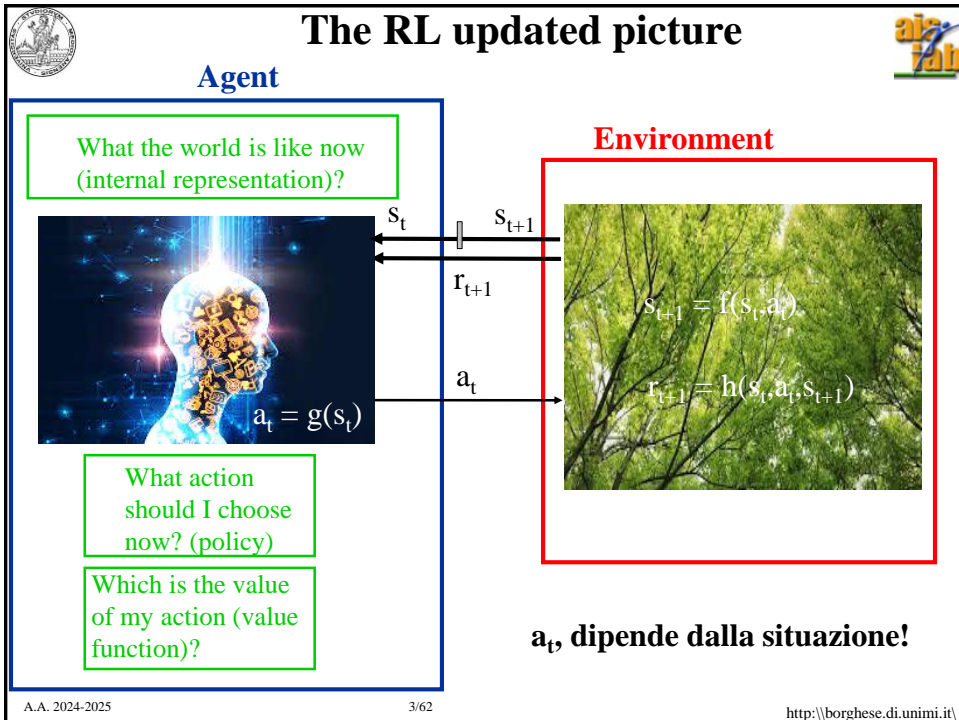
Iterative policy evaluation

Miglioramento della policy

A.A. 2024-2025

2/62

<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) Apprendo la sua Value function ($Q^\pi(s,a)$)**
- 3) Miglioro la policy.

Itero i passi 2 e 3 fino a quando non raggiungo l'ottimo.

A.A. 2024-2025

4/62

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



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 $Q^\pi(\cdot)$.
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 $Q^\pi(\cdot)$.
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 $Q^\pi(\cdot)$.



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.

Azioni ammissibili (policy):

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

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

Goal: collezionare il maggior numero di lattine.



Esempio di calcolo della Value function



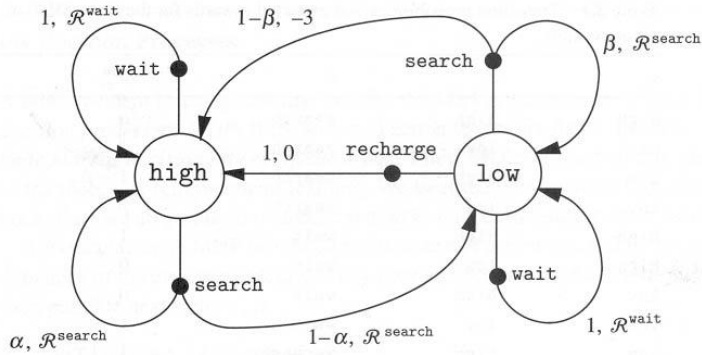
Policy deterministica

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

Value function

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

$\alpha = \Pr(s_{t+1} = \text{High} | s_t = \text{High}, a_t = \text{Search}) = 0.4$
 $\beta = \Pr(s_{t+1} = \text{Low} | s_t = \text{Low}, a_t = \text{Search}) = 0.1$
 $\gamma = 0.8, R^{\text{search}} = 3, R^{\text{wait}} = 1, R^{\text{dead}} = -3, R^{\text{auto}} = 0$



A.A. 2024-2025

7/62

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



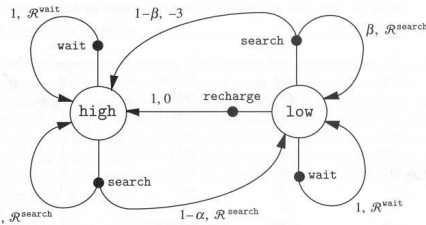
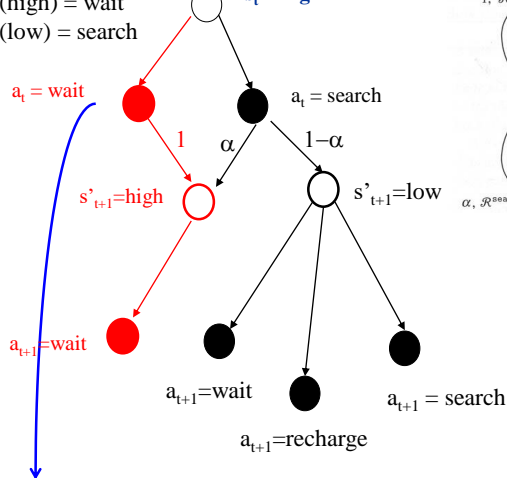
Analisi a un passo al tempo t - $s_t = \text{high}$



Policy deterministica

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

$s_t = \text{high}$



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

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

A.A. 2024-2025

8/62

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

Analisi a un passo al tempo t - $s_t = \text{high}$

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

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

$s_t = \text{high}$

$Q^\pi(\text{high}, \text{wait}) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} = R^{\text{wait}} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} = R^{\text{wait}} + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k+1} = R^{\text{wait}} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2}$

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

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

A.A. 2024-2025 9/62 http://borghese.di.unimi.it/

Analisi a un passo al tempo t - $s_t = \text{low}$

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

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


$s_t = \text{low}$

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


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

2 cammini possibili!!

A.A. 2024-2025 10/62 http://borghese.di.unimi.it/



Policy deterministica - $s_t = \text{low}$



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

$s = \text{High} - a = \text{Wait};$
 $s = \text{Low} - a = \text{Search};$


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

$$Q^\pi(\text{low}, \text{search}) = 0.1 \times [3 + 0.8 \times Q^\pi(\text{low}, \text{search})] + 0.9 \times [-3 + 0.8 \times Q^\pi(\text{high}, \text{wait})]$$


A.A. 2024-2025

11/62

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



Analisi ad un passo dal tempo t



Policy deterministica

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

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

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

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/>



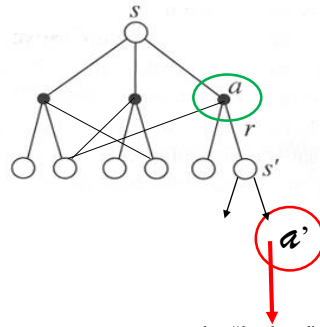
Calcolo ricorsivo della Value function



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

$$Q^\pi(s_{t+1}, a_{t+1}) = E_\pi\{R_t | s_{t+1} = s', a_{t+1} = a'\} = \sum_{k=0}^{\infty} \gamma^k r_{t+k+2}$$

Relazione tra $Q^\pi(s, a)$ e $Q^\pi(s_{t+1}, a_{t+1})$?



A.A. 2024-2025

13/62

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



Calcolo ricorsivo della Value function



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

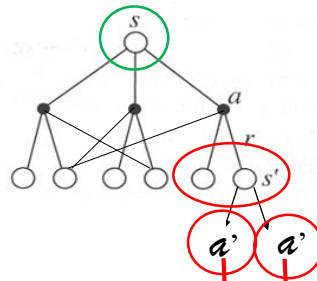
Isolo il reward ad un passo nella serie dei reward.

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

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

Io termine
(a un passo)

Io termine
(passi futuri)



A.A. 2024-2025

14/62

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



$Q^\pi(s_t, a_t)$: primo termine

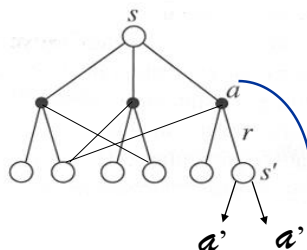


$$P_{s \rightarrow s' | a} \triangleq \Pr(s_{t+1} = s' | s_t = s, a_t = a)$$

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

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

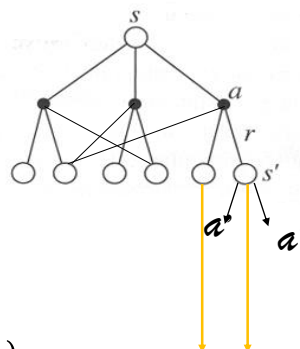


$Q^\pi(s_t, a_t)$: secondo termine




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

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




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

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



Putting all together



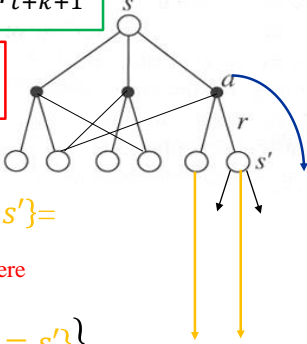
$$Q^\pi(s_t, a_t) = E_\pi\{R_t | s_t = s, a_t = a\} = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

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

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


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

Io termine (a un passo) Io termine (passi futuri)




Not yet there

A.A. 2024-2025
17/62
<http://borghese.di.unimi.it/>



Formulazione ricorsiva - policy deterministica



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

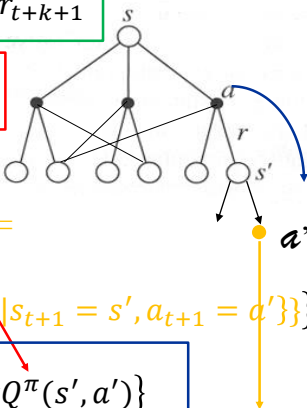
$$Q^\pi(s_{t+1}, a_{t+1}) = E_\pi\{R_t | s_{t+1} = s', a_{t+1} = a'\}$$

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

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

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

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



A.A. 2024-2025
18/62
<http://borghese.di.unimi.it/>



Osservazioni

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

Devo considerare i reward a un passo che portano da s a tutti gli stati prossimi s' che possono venire visitati.

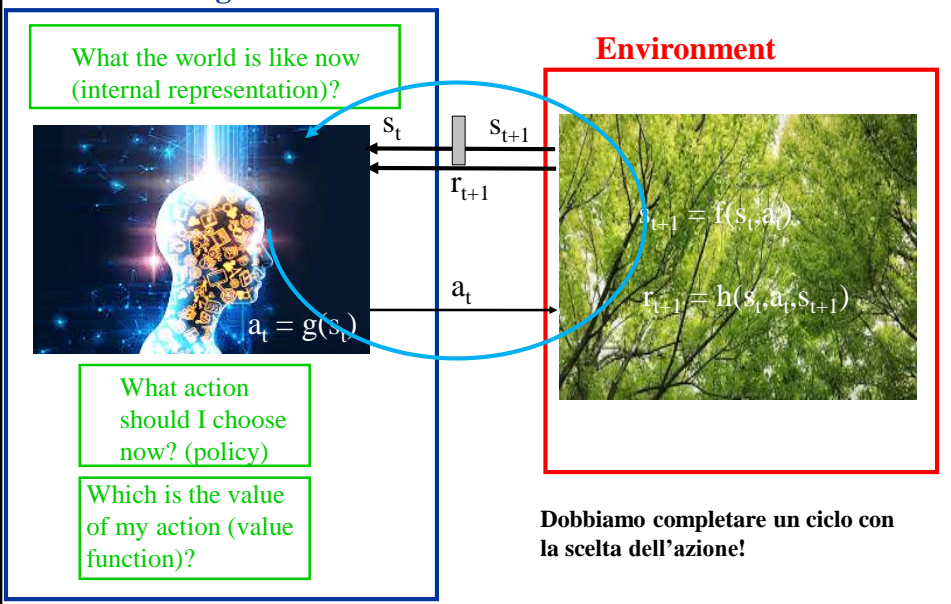
A partire da ogni s', devo considerare il reward a lungo termine che si può accumulare nell'interazione con l'ambiente, $Q^\pi(s', a')$.



Un ciclo di interazione

Agent

Environment





Sommario



Le equazioni di Bellman per policy deterministica

Le equazioni Bellman per policy stocastica

Iterative policy evaluation

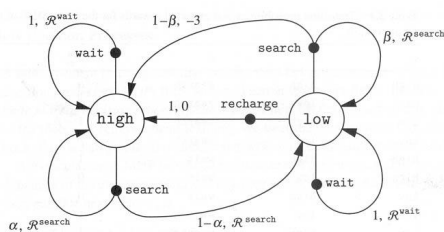
Miglioramento della policy



Valutazione policy stocastica

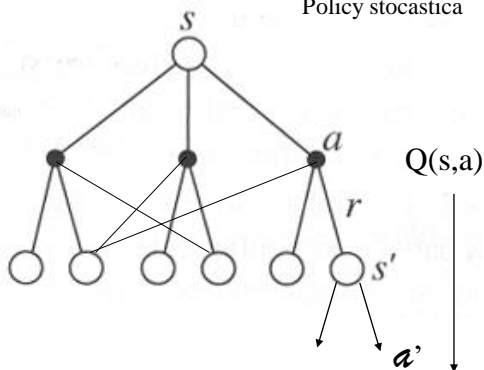


Nel valutare $Q(s,a)$ dobbiamo valutare tutti i cammini che partono da ogni s' .



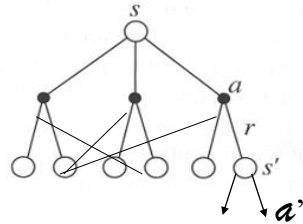
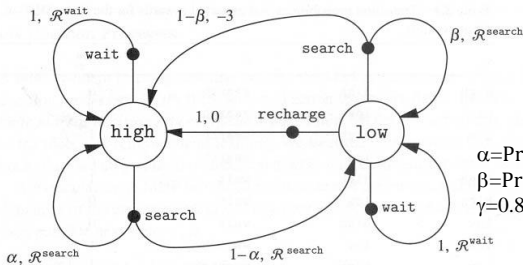
$\pi(s,a)$ stocastica

Policy stocastica





Policy stocastica



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

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

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

Deterministica (1 azione scelta in s'):

$$Q(\text{high}, \text{wait}) = 1 \times \{ \mathcal{R}^{\text{wait}} + \gamma [Q(\text{high}, \text{wait})] \}$$

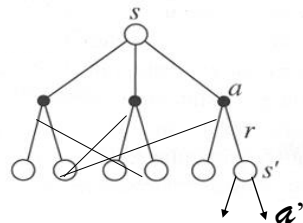
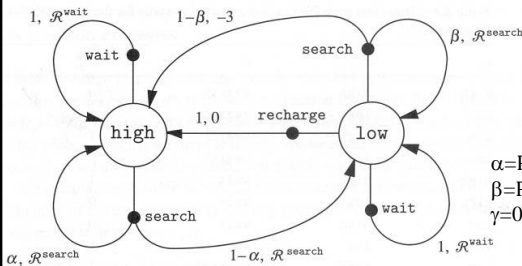
Stocastica (più azioni scelte in s'):

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

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



Policy stocastica



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

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

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


Deterministica (1 azione scelta in s'):

$$Q(\text{high}, \text{search}) = \Pr(s_{t+1} = \text{High} | s_t = \text{High}, a_t = \text{Search}) \times \{ \mathcal{R}^{\text{search}} + \gamma [Q(\text{high}, \text{search})] \} + \Pr(s_{t+1} = \text{Low} | s_t = \text{High}, a_t = \text{Search}) \times \{ \mathcal{R}^{\text{search}} + \gamma [Q(\text{low}, \text{search})] \}$$

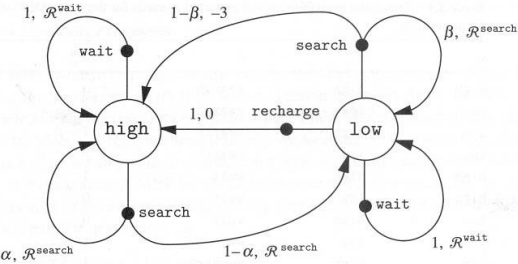
Stocastica (più azioni scelte in s'):

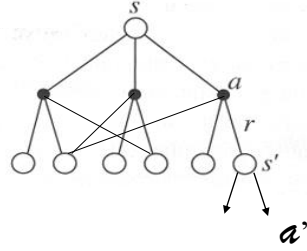
$$Q(\text{high}, \text{search}) = \Pr(s_{t+1} = \text{High} | s_t = \text{High}, a_t = \text{Search}) \times \{ \mathcal{R}^{\text{search}} + \gamma [(\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait}))] \} + \Pr(s_{t+1} = \text{Low} | s_t = \text{High}, a_t = \text{Search}) \times \{ \mathcal{R}^{\text{search}} + \gamma [\Pr(a' = \text{search} | \text{low}) Q(\text{low}, \text{search}) + \Pr(a' = \text{wait} | \text{low}) Q(\text{low}, \text{wait}) + \Pr(a' = \text{recharge} | \text{low}) Q(\text{low}, \text{rech})] \}$$

$$Q(\text{high}, \text{search}) = 0.4 \times \{ 3 + 0.8 [(\Pr(a' = \text{search} | \text{high}) Q(\text{high}, \text{search}) + \Pr(a' = \text{wait} | \text{high}) Q(\text{high}, \text{wait}))] \} + 0.6 \times \{ 3 + 0.8 [\Pr(a' = \text{search} | \text{low}) Q(\text{low}, \text{search}) + \Pr(a' = \text{wait} | \text{low}) Q(\text{low}, \text{wait}) + \Pr(a' = \text{recharge} | \text{low}) Q(\text{low}, \text{rech})] \}$$



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


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

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


A.A. 2024-2025

5 equazioni in 5 incognite

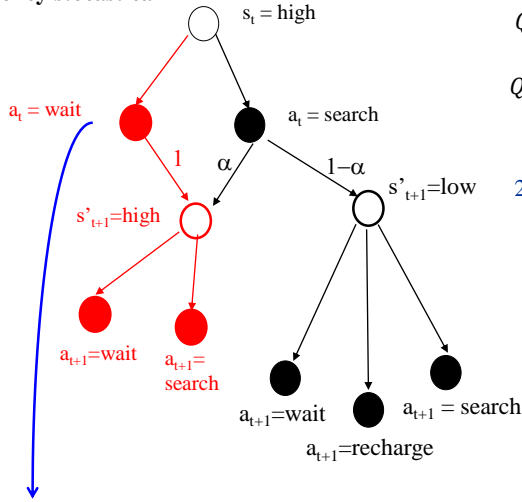
<http://borgese.di.unimi.it/>



Analisi a un passo al tempo t



Policy stocastica

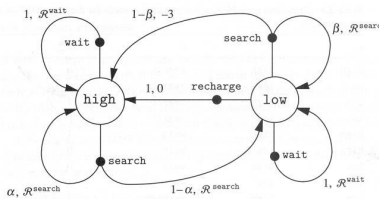


$$Q^\pi(s_t, a_t) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

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

2 cammini possibili in s_t !!


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




A.A. 2024-2025

26/62

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



Analisi a un passo al tempo t



Policy stocastica (uniforme)

Wait e search equiprobabili

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

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

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. 2024-2025

27/62

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

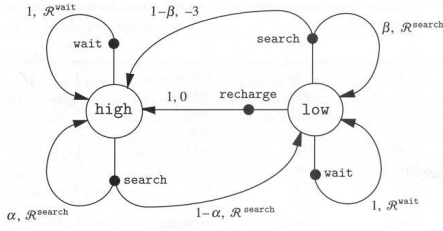


Analisi a un passo al tempo t



Policy stocastica

5 cammini possibili!!



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


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

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

A.A. 2024-2025


28/62

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



Analisi a un passo al tempo t

Policy stocastica (equiprobabile)



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

5 cammini possibili!!

Recharge, wait e search equiprobabili

Wait e search equiprobabili


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

A.A. 2024-2025


29/62

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



Analisi a un passo al 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. 2024-2025

30/62

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



Formulazione ricorsiva - policy stocastica



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

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

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

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

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

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

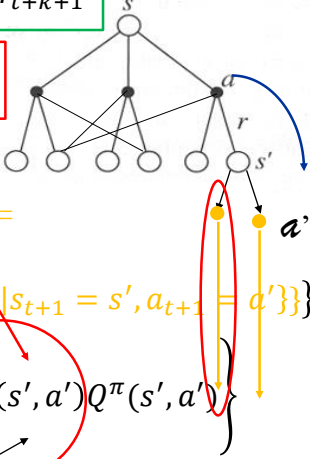
Io termine
(a un passo)


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

A.A. 2024-2025


31/62

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





Sommaro



- Le equazioni di Bellman per policy deterministica
- Le equazioni Bellman per policy stocastica
- Iterative policy evaluation
- Miglioramento della policy

A.A. 2024-2025

32/62

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



Fondamenti del metodo



- Supponiamo di essere all'istante t . In questo istante t , siamo in s_t e da s_t si può passare a un certo insieme di stati: $\{s'_{t+1}\}$.
- Analizziamo un solo passo: cosa succede nella transizione da t a $t+1$.
- Migliorare la stima della nostra Value Function ad ogni iterazione.

A.A. 2024-2025

33/62

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

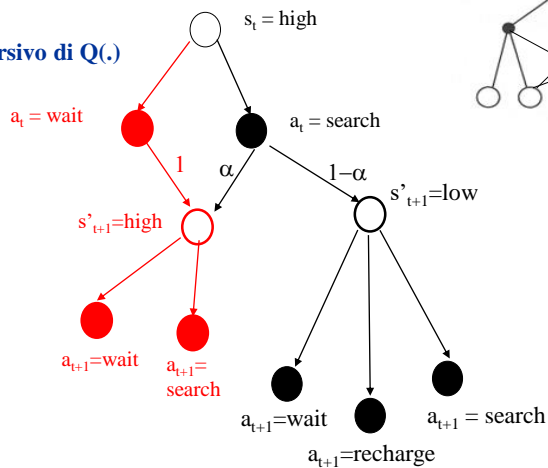


Equazioni di Bellman



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

Calcolo ricorsivo di $Q(\cdot)$



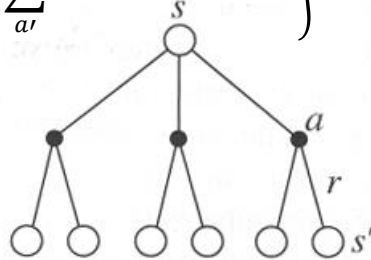
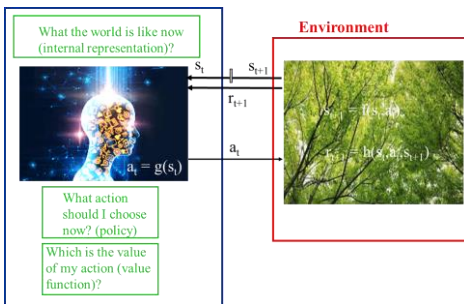
$$Q^\pi(\text{high}, \text{wait}) = R^{\text{wait}} + 0.5 \gamma Q^\pi(\text{high}, \text{wait}) + 0.5 \gamma Q^\pi(\text{high}, \text{search})$$



Osservazioni

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

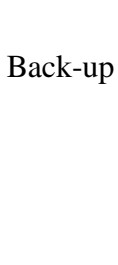
Calcolo ricorsivo di Q(.)



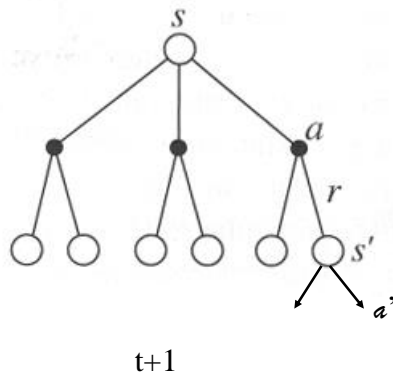
Passo da t a t+1 poi guardo backwards in time



Tecnica full-back



$\pi(s,a)$ fissata



Conosciamo $Q^\pi(s_t, a_t) \forall s_t, a_t$ 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^π per $\{s, a\}$: $Q^\pi(s_t, a_t)$ congruente con:

$$Q^\pi(s_t, a_t) \text{ ed } r_{t+1}$$

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

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

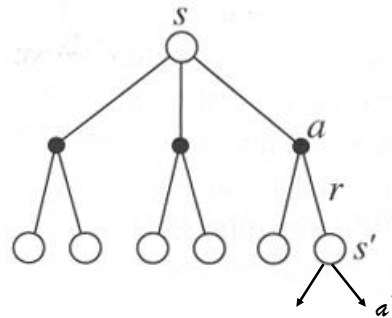
π fissata



Tecnica full-back

Back-up

$\pi(s,a)$ fissata



t+1

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

Conosciamo $Q^\pi(s_t, a_t) \forall s_t, a_t$ 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^π per $\{s, a\}$: $Q^\pi(s_t, a_t)$ congruente con:

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

π fissata

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

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

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



Algoritmo per "iterative policy evaluation", versione batch



Partiamo da una politica $\pi(s,a)$ data, supponiamo deterministica.

Definiamo una soglia di convergenza τ

Inizializziamo $Q(s,a) = 0 \forall s, \forall a$, compreso gli stati finali.

Repeat

```

{  Δ = 0;
  for s = 1 : NS
    {  a = policy(s);
      {  Temp_Q(s,a) = 0;
        for s_next = 1 : NS
          {  Pr_s_next = NextState(s,a);
            reward = ComputeReward(s,a,s_next);
            a_next = policy(s_next);
            Temp_Q(s,a) = Pr_s_next * (reward + γQ(s_next,a_next));
          }
        }
      }
    }
  for s=1:NS;
  {
    a = policy(s);
    if ( | Temp_Q(s,a) - Q(s,a) | > Δ )
    {
      Δ = | Temp_Q(s,a) - Q(s,a) |;
    }
    Q(s,a) = Temp_Q(s,a);
  }
} Until (Δ < τ);

```

1 step Forwards
Pass for all states

1 step Backwards
Pass for all states
(full backup)



Interpretazione dell'update (batch o trial)



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

Al termine dell'aggiornamento dei $Q^\pi(s,a)$ per tutti gli stati,
 $Q^\pi(s,a) = Q^\pi_{\text{new}}(s,a)$. **Aggiornamento batch.**

In alternativa, utilizzerò in parte già il nuovo valore di
 $Q^\pi(s,a)$ all'interno dell'equazione di aggiornamento.
Aggiornamento per trial.

Entrambe le modalità di aggiornamento convergono.



Algoritmo per "iterative policy evaluation", versione per trial



Partiamo da una politica $\pi(s,a)$ data, deterministica.

Definiamo una soglia **relativa** di convergenza τ .

Inizializziamo $V(s) = 0 \forall s$, compreso gli stati finali.

Repeat

```

{
  Δ = 0;
  for s = 1 : NS
    a = policy(s);
    Value = Q(s,a);
    for s_next = 1 : NS
      Pr_s_next = NextState(s,a);
      reward = ComputeReward(s,a,s_next);
      a_next = policy(s_next);
      Q(s,a) = Pr_s_next * (reward + γQ(s_next,a_next));
      Δ = max(Δ, (| Value - Q(s,a) |));
    }
  }
} Until (Δ < τ);
  
```

1 step Forwards
Pass for all states

// Until end of States – End of an episode



Problematiche legate al calcolo di $V(s)$: problema di policy evaluation



3 assunzioni:

- 1) Conoscenza della dinamica dell'ambiente: $P(s \rightarrow s' | a)$
- 2) Conoscenza della policy (eventualmente stocastica), $\pi(s, a)$
- 3) Potenza di calcolo sufficiente
- 4) Proprietà Markoviane dell'ambiente (definizione di uno stato).

Le equazioni contengono dei termini statistici (valori attesi).

Soluzione di un sistema lineare in N incognite (numero di stati).

Come mai posso determinare la Value function per la policy $\pi(\cdot)$, se questa si basa sul reward che riceverò negli istanti futuri?

C'è poca interazione con l'ambiente e molta simulazione (cf. metodi Montecarlo).



Riassunto



Posso determinare la Value function in modo ricorsivo. Per ogni stato, sarà funzione dell'output dell'ambiente in quell'istante (attraverso la funzione stato prossimo ed il reward istantaneo) e della policy scelta in quell'istante e dei reward a lungo termine attesi negli stati in cui l'ambiente mi porta.

Per scegliere la policy devo esaminare il reward a lungo termine che mi si prospetta nello stato in cui mi trovo e scegliere l'azione che lo massimizza.



Sommario



Le equazioni di Bellman per policy deterministica

Le equazioni Bellman per policy stocastica

Iterative policy evaluation

Miglioramento della policy



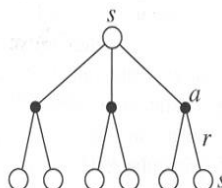
Meccanismo di apprendimento nel RL



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.**





Miglioramento della policy



Tutti gli stati sono valutati in funzione di una policy data.

Condizioni di funzionamento dell'agente:

- Policy **deterministica**: $a = \pi(s)$.
- Ambiente **stocastico**.

Cosa succede se cambiamo la policy per un certo stato s_m ? $a_{new} \neq \pi(s_m)$.
Cosa viene influenzato?

Scelgo a_{new} in s_m , visiterò una certa sequenza di stati, per questi stati seguirò la policy precedente per $s \neq s_m$. Cosa viene influenzato?

Come faccio a valutare se miglioro la policy o no?



Effetto del cambiamento della policy



Cambia, a, cambiano i possibili stati successivi ad s_m , $\{s_{t+k}\}$, ed il reward a lungo termine ($V^\pi(s_{t+1}) = Q^\pi(s_{t+1}, a_{t+1})$) per policy deterministica):

$$Q^\pi(s_m, a_{new}) = E_\pi \{ r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s_m, a_t = a_{new} \neq \pi(s_m) \} =$$

$$\sum_{s'} P_{s_m \rightarrow s'}^{a_{new}} [R_{s_m \rightarrow s'}^{a_{new}} + \gamma V^\pi(s')]]$$

?

$$Q^\pi(s_m, a_{new}) \geq Q^\pi(s_m, a = \pi(s_m)) \quad \forall s, a?$$

Se il reward fosse migliore con a_{new} , sceglierò sempre a_{new} in s_m .

Il reward a lungo termine può essere maggiore (minore) solamente se aumenta (diminuisce) il reward totale "visto" a un passo (reward del passo + reward successivo).



Enunciato del teorema del miglioramento della policy

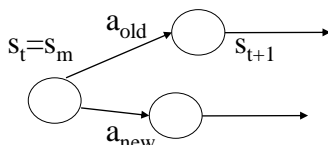


$$Q^{\pi}(s, a) = \sum_k P_{s \rightarrow s_k | a} [R_{s \rightarrow s_k | a} + \gamma V^{\pi}(s_k)]$$

Ipotesi: π and π' deterministic policies
 $Q^{\pi}(s_m, \pi'(s_m)) \geq V^{\pi}(s_m)$

$$Q^{\pi}(s, a_{new} = \pi'(s_m)) = \sum_k P_{s_m \rightarrow s_k | a_{new}} [R_{s_m \rightarrow s_k | a_{new}} + \gamma V^{\pi}(s_k)]$$

Tesi: π' è meglio di π . Cioè: $V^{\pi'}(s, a(s)) \geq V^{\pi}(s, a(s)) \forall s$.
 $Q^{\pi'}(s, a_{new}) \geq Q^{\pi}(s, a_{old})$



A.A. 2024-2025

47/62

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



Dimostrazione del teorema del miglioramento della policy



Analizziamo la seguente condizione:

$\pi' = \pi \forall s$ tranne che per s_m per il quale si applica l'azione:
 $a_{new} = \pi'(s_m)$

Risulta che il reward a lungo termine è maggiore per $a_{new} = \pi'(s)$.

$$V^{\pi'}(s) = Q^{\pi'}(s, a_{new} = \pi'(s)) \geq Q^{\pi}(s, a = \pi(s)) = V^{\pi}(s)$$

Tesi: π' è meglio di π . Cioè: $V^{\pi'}(s) \geq V^{\pi}(s) \forall s$ (ed in particolare per gli altri stati s)

A.A. 2024-2025

48/62

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



Dimostrazione del teorema del miglioramento della policy



Hp: $Q^\pi(s, \pi'(s)) \geq V^\pi(s) \quad \forall s \quad \pi'(s, a)$ è migliore per almeno uno stato

$$V^\pi(s) \leq Q^\pi(s, \pi'(s))$$

$$= E_{\pi'}\{r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s\}$$

$$\leq E_{\pi'}\{r_{t+1} + \gamma Q^\pi(s_{t+1}, \pi'(s_{t+1})) \mid s_t = s\}$$

$$\leq E_{\pi'}\{r_{t+1} + \gamma E_{\pi'}(r_{t+2} + \gamma V^\pi(s_{t+2})) \mid s_t = s\}$$

$$= E_{\pi'}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 V^\pi(s_{t+2}) \mid s_t = s\}$$

Sostituisco ancora $Q^{\pi^*}(\cdot)$

$$\leq E_{\pi'}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s_t = s\}$$

Th: $V^\pi(s) \leq V^{\pi'}(s)$



Osservazioni



$$s = s_m \quad Q^\pi(s_m, \pi'(s)) \geq Q^\pi(s_m, \pi(s))$$

$$s \neq s_m \quad Q^\pi(s, a) = E_{\pi'}\{r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s\}$$

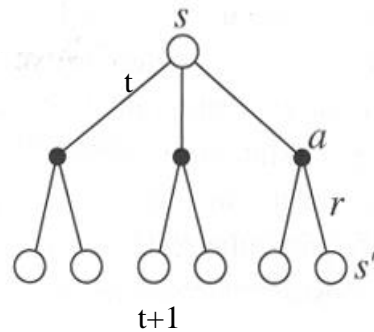
$$= E_{\pi'}\{r_{t+1} + \gamma Q^\pi(s_{t+1}, \pi(s_{t+1})) \mid s_t = s\}$$

Se $s_{t+k} = s_m$ miglioro la $Q(s, a)$.

Se nessun $s_{t+k} = s_m$, Non varia la $Q(s, a)$.



Visione grafica del miglioramento



Ogni volta che sono in uno stato, s , scelgo un'azione che migliora il reward a lungo termine ottenuto da quell'istante/stato in poi.

Per gli altri stati, il reward a lungo termine non viene modificato ogni volta che l'albero uscente da s passa per s .



Ottimizzazione policy



Per ogni stato scelgo le azioni secondo la policy: $\pi(s,a)$.

Posso ordinare la Value function $Q(s,a)$ in ordine decrescente, in funzione delle azioni scelte in s (policy).


Si definisce una policy, π_1 , migliore di un'altra, π_2 , se e solo se:

$$Q^{\pi_1}(s,a(s)) \geq Q^{\pi_2}(s,a(s)) \quad \forall s.$$


In particolare si definisce una politica ottima, π^* , se e solo se:

$$Q^*(s,a(s)) \geq V^{\pi}(s,a(s)) \quad \forall s$$

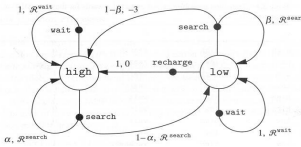
$$Q^*(s,a(s)) \geq Q^{\pi}(s,a(s)) \quad \forall [s,a]$$



Q(s, a) - Osservazioni



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

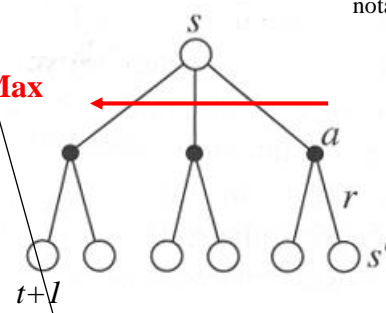


1, q_{wait}
1- β , -3
 β , $q_{recharge}$
1, 0
 α , $q_{recharge}$
1- α , $q_{recharge}$
1, q_{wait}

Policy nota


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

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




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

A.A. 2024-2025
53/62
<http://borghese.di.unimi.it/>



Calcolo ricorsivo della Value function ottima: confronti



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

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

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

Politica greedy: scelgo l'azione ottimale.
Ha senso per il robot raccogli-lattine?

A.A. 2024-2025
54/62
<http://borghese.di.unimi.it/>



Policy iteration

Iterazione tra:

- Calcolo iterativo della Value function (iterative policy evaluation)
- Miglioramento della policy (policy improvement)

$$\begin{array}{ccccccccccc} \pi_0 & \rightarrow & V^{\pi_0} & \rightarrow & \pi_1 & \rightarrow & V^{\pi_1} & \rightarrow & \pi_2 & \rightarrow & V^{\pi_2} & \rightarrow & \dots \\ & & & & \rightarrow & & \rightarrow & & & & & & \end{array}$$

Converge velocemente ad una buona politica
(cf. Software Sommaruga)



Algoritmo

Inizialization

$$Q(s,a) = 0;$$

$$\pi(s,a) = \text{random (e.g. equiprobabile);}$$

Repeat

 Policy evaluation.

 Policy improvement.

until policy_stable



Algoritmo per "iterative policy evaluation", versione per trial



Partiamo da una politica $\pi(s,a)$ data, deterministica.
Definiamo una soglia **relativa** di convergenza τ .
Inizializziamo $V(s) = 0 \forall s$, compreso gli stati finali.

```
Repeat
{
   $\Delta = 0$ ;
  for s = 1 : NS //  $\forall s, \neq TS$  1 step Forwards
  {
    a = policy(s); Pass for all states
    Value = Q(s,a);
    for s_next = 1 : NS // for all next states
    {
      Pr_s_next = NextState(s,a); // compute the probability that  $s_{t+1} = s'$ 
      reward = ComputeReward(s,a,s_next); // Compute average 1 step reward
      a_next = policy(s_next); // Next action
       $Q(s,a) = Pr\_s\_next * (reward + \gamma Q(s\_next, a\_next))$ ;
       $\Delta = \max(\Delta, (| Value - Q(s,a) |))$ ;
    }
  }
} // Until end of States – End of an episode

} Until ( $\Delta < \tau$ );
```



Policy improvement



```
policy_stable = true;
for s = 1:NS // in alternativa, scelgo uno stato
  a_old =  $\pi(s)$ ;
  a_new =  $\arg \max_{a'} \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s,s',a} + \gamma \pi(s', a') \sum_{a'} Q^\pi(s', a') \right\}$ 
  if (a_new  $\neq$  a_old)
    policy_stable = false;
end;
```

Operazione di max hard o «soft» -> policy ϵ -greedy, pursuit, ...



Iterative policy evaluation sulla value function $V(s)$



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

Converge al limite a $Q^{\pi}(s, a)$. Come facciamo a troncare?



Value iteration



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

Invece di considerare una policy stocastica, consideriamo l'azione migliore:

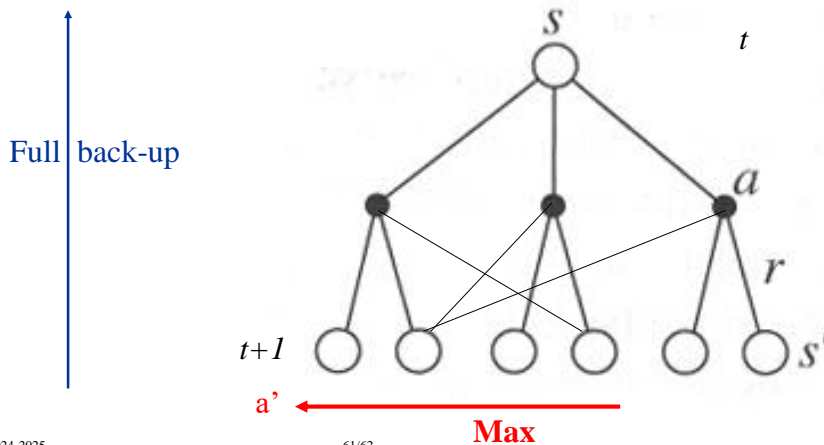
$$Q_{k+1}^{\pi}(s_t, a_t) = \max_{a'} \sum_{s'} P_{s \rightarrow s' | a} \left\{ R_{s, s', a} + \gamma \pi(s', a') \sum_{a'} Q_k^{\pi}(s', a') \right\}$$

$\forall s$



Visualizzazione grafica

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



A.A. 2024-2025

61/62

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



Sommario

- Le equazioni di Bellman per policy deterministica
- Le equazioni Bellman per policy stocastica
- Iterative policy evaluation
- Miglioramento della policy

A.A. 2024-2025

62/62

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