

Sistemi Intelligenti Reinforcement Learning: Equazioni di Bellman

Alberto Borghese

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



A.A. 2006-2007

1/26

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



Sommario



Le equazioni di Bellman

Value function ottima

A.A. 2006-2007

2/26

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



Il modello markoviano

Il comportamento dell'ambiente è definito dallo stato: $S = \{s_j\}$
Per ogni stato l'agente sceglie un'azione: $a = a(s)$ $A = \{a_k\}$
Policy di un agente: $\pi(s, a)$ è quanto dobbiamo definire.

L'ambiente ha una evoluzione stocastica rappresentata da un MDP:

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

Inoltre, ad ogni istante fornisce un reward immediato associato alla transizione, stimato all'istante t come:

$$R_{s_t=s \rightarrow s_{t+1}=s' | a_t=a} = E\{r_{t+1} = r' | s_t = s, a_t = a, s_{t+1} = s'\}$$

$$\forall s \in S; \forall a \in A$$



La value function

Nulla è detto sulla policy: dato uno stato, quale azione scegliere? In quale stato mi sposto?

Vogliamo costruire agenti lungimiranti.

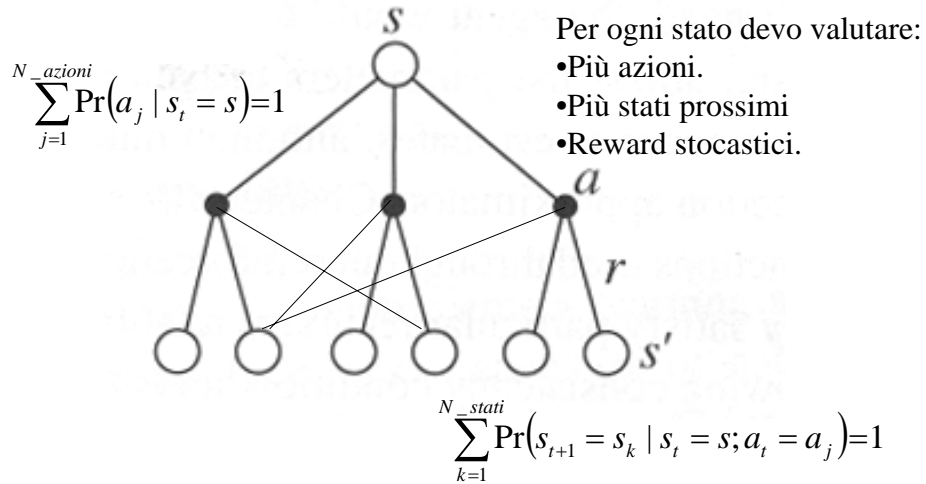
State-Value function:

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

Massimizzo la ricompensa a lungo termine, $V(\cdot)$. Dipende dalla policy:



Value function e modelli markoviani



A.A. 2006-2007

5/26

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



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.

Policy:

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

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

Task: collezionare il maggior numero di lattine.

A.A. 2006-2007

6/26

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



Funzionamento del Robot



Funzione Stato prossimo:

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

Se il livello di energia è alto ($s_t = \text{alto}$):

se scelgo Wait - $s_{t+1} = \text{alto}$.

se scelgo Search, s_{t+1} avrà una certa probabilità di diventare low.

$$P_{\text{high} \rightarrow \text{low} | \text{Search}} = \Pr\{s_{t+1} = \text{low} | s_t = \text{high}, a_t = \text{Search}\} = \alpha$$

Se il livello di energia è basso ($s_t = \text{basso}$):

se scelgo Wait - $s_{t+1} = \text{basso}$.

se scelgo Recharge - $s_{t+1} = \text{alto}$.

se scelgo Search, s_{t+1} avrà una certa probabilità di fermarsi.

$$P_{\text{low} \rightarrow \text{low} | \text{Search}} = \Pr\{s_{t+1} = \text{low} | s_t = \text{low}, a_t = \text{Search}\} = \beta$$

A.A. 2006-2007

7/26

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



Reward del Robot



Funzione Reward:

$$R_{s \rightarrow s'|a} = E\{r_{t+1} = r' | s_t = s, a_t = a, s_{t+1} = s'\}$$

R^{search} reward se il robot sta cercando.

R^{wait} reward se il robot sta cercando.

-3 se occorre portarlo a ricaricarsi.

0 se il robot va autonomamente a ricaricarsi.

$$R^{\text{search}} > R^{\text{wait}}$$

A.A. 2006-2007

8/26

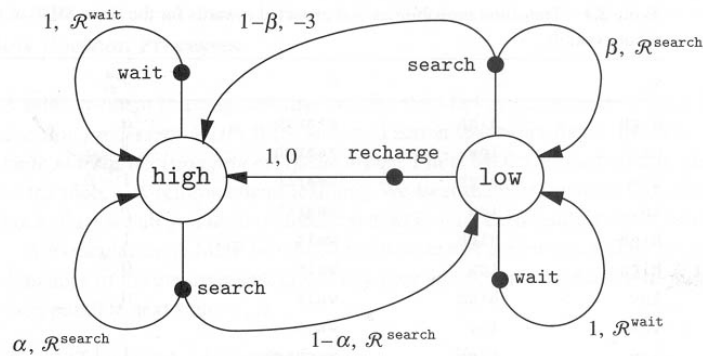
<http://homes.dsi.unimi.it/~borghese/>



Esempio di calcolo della funzione valore



$V(\text{high}) = ?$ $a(\text{high}) = ?$
 $V(\text{low}) = ?$ $A(\text{low}) = ?$



A.A. 2006-2007

9/26

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



Calcolo ricorsivo della Value function



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

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

Relazione?

$$V^\pi(s) = E_\pi \left\{ r_{t+1} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s \right\} =$$

$$V^\pi(s) = E_\pi \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_t = s \right\}$$

A.A. 2006-2007

10/26

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

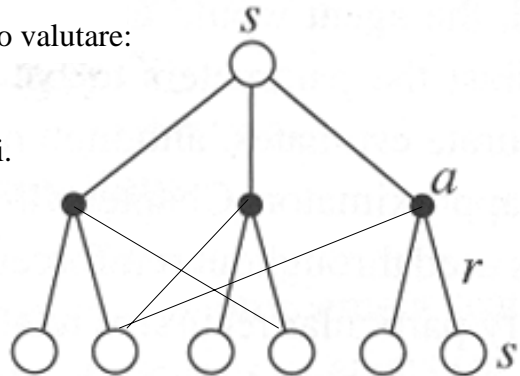


V(s) : primo termine

$$E_{\pi} \{ r_{t+1} \mid s_t = s \} = \left[\sum_{a_j} \pi(a_j, s) \right] \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j}]$$

Per ogni stato devo valutare:

- Più azioni.
- Più stati prossimi
- Reward stocastici.



A.A. 2006-2007

11/26

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

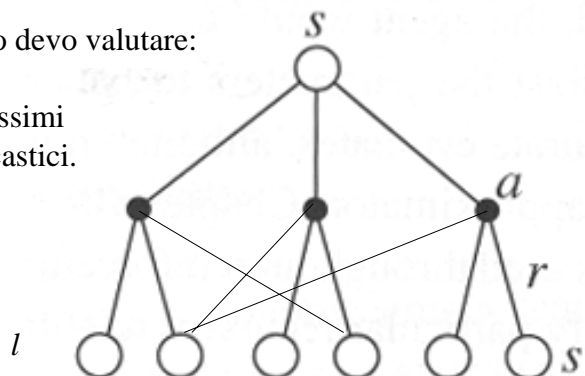


V(s) : secondo termine

$$E_{\pi} \left\{ \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_t = s \right\} = \gamma \sum_l E_{\pi} \{ R_{t+1} \mid s_{t+1} = s'_l \} \Pr(s_{t+1} = s'_l \mid s_t = s)$$

Per ogni stato devo valutare:

- Più azioni.
- Più stati prossimi
- Reward stocastici.



A.A. 2006-2007

12/26

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



Calcolo ricorsivo della Value function



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

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

Legame?

Policy Next-state

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

$$V^\pi(s) = \sum_{a_j} \pi(a_j, s) \sum_{s'} P_{s \rightarrow s' | a_j} R_{s \rightarrow s' | a_j} + E_\pi \left\{ \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_t = s \right\}$$

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

Bellman's equation

A.A. 2006-2007

13/26

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

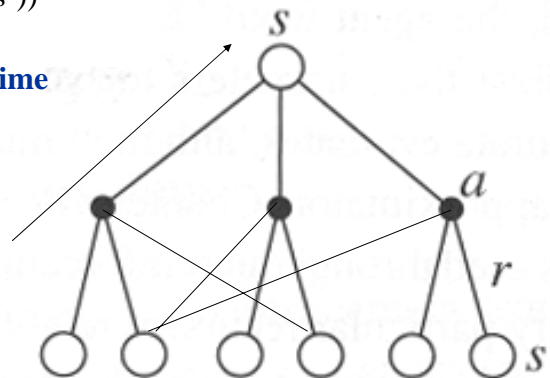


Osservazioni



$$V^\pi(s) = \text{funz}(V^\pi(s'))$$

Backwards in time



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

A.A. 2006-2007

14/26

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



Reinforcement Learning Problem



Given: Repeatedly...

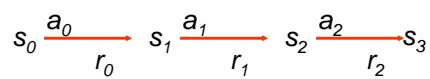
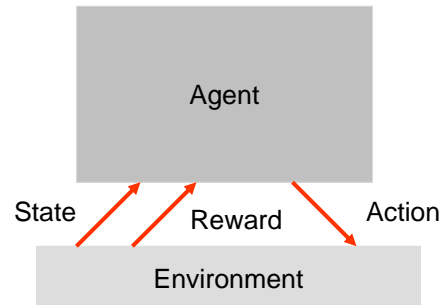
- Executed action
- Observed state
- Observed reward

Learn action policy $\pi: S \rightarrow A$

- ◆ Maximizes life reward
 $r_0 + \gamma r_1 + \gamma^2 r_2 \dots$
from any start state.
- ◆ Discount: $0 < \gamma < 1$

Note:

- Unsupervised learning
- Delayed reward



Goal: **Learn** to choose actions that maximize life reward

$$r_0 + \gamma r_1 + \gamma^2 r_2 \dots$$



Sommario



Le equazioni di Bellman

Value function ottima



Ottimizzazione Value function e policy



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

Posso ordinare la Value function $V(s)$ in funzione delle azioni scelte in s (policy).

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

$$V^{\pi_1}(s) \geq V^{\pi_2}(s) \quad \forall s.$$

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

$$V^*(s) \geq V^\pi(s) \quad \forall s$$

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



Calcolo ricorsivo della Value function ottima



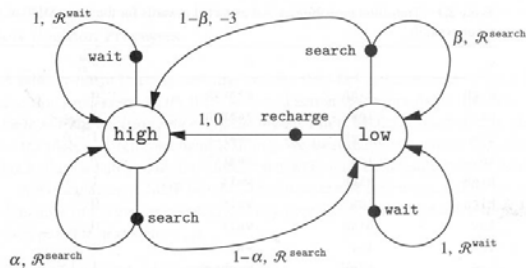
$$V^\pi(s) = \left\{ \sum_{a_j} \pi(a_j, s) \sum_{s_l'} P_{s \rightarrow s_l' | a_j} \left[R_{s \rightarrow s_l' | a_j} + \gamma V^\pi(s_l') \right] \right\}$$

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

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



Esempio di calcolo della funzione valore



$\alpha=0.6, \beta=0.1, \gamma=0.8, R^{\text{search}}=3, R^{\text{wait}}=1$
 $\text{Pr}(a, \text{state}=\text{high}) = [0.4, 0.6, 0]$
 $\text{Pr}(a, \text{state} = \text{low}) = [0.4, 0.5, 0.1]$

$$V_h = 0.4 \times 1 \times [1 + 0.8V_h] + 0.6 \times 0.4 \times [3 + 0.8V_h] + 0.6 \times 0.6 \times [3 + 0.8V_l]$$

Wait *Search -> high* *Search -> low*

$$V_l = 0.4 \times 1 \times [1 + 0.8V_l] + 0.5 \times 0.1 \times [3 + 0.8V_l] + 0.6 \times 0.9 \times [-3 + 0.8V_h] + 0.1 \times 1 \times [0 + 0.8V_h]$$

Wait *Search -> Low* *Recharge* *Search -> High*

Sistema lineare di 2 equazioni nelle 2 incognite: V_h e V_l
Come calcolo V_h^* e V_l^* ?

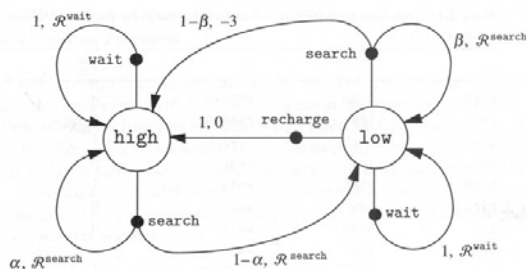
A.A. 2006-2007

19/26

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



Esempio di calcolo della policy ottima



$V(\text{high})$:

$$\text{state} = \text{high}, \text{action} = \text{wait} \rightarrow Q(\text{high}, \text{wait}) = 1[R_{\text{wait}} + \gamma V^*(\text{high})]$$

$$\text{state} = \text{high}, \text{action} = \text{search} \rightarrow Q(\text{high}, \text{search}) = \alpha[R^{\text{search}} + \gamma V^*(\text{high})] + (1 - \alpha)[R^{\text{search}} + \gamma V^*(\text{low})]$$

$V(\text{low})$:

$$\text{state} = \text{low}, \text{action} = \text{wait} \rightarrow Q(\text{low}, \text{wait}) = [R_{\text{wait}} + \gamma V^*(\text{low})]$$

$$\text{state} = \text{low}, \text{action} = \text{search} \rightarrow Q(\text{low}, \text{search}) = \beta[R^{\text{search}} + \gamma V^*(\text{low})] + (1 - \beta)[-3 + \gamma V^*(\text{high})]$$

$$\text{state} = \text{low}, \text{action} = \text{rechar} \rightarrow Q(\text{low}, \text{rechar}) = \gamma V^*(h)$$

A.A. 2006-2007

20/26

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



Utilizzo di $V(s)$



Sistema di N equazioni non-lineari in N incognite (le entry della tabella della policy).

Politica greedy rispetto alla Value function.

Questa politica greedy (ad un passo) produce una politica ottima globalmente.

Vengono valutate le conseguenze a breve termine delle azioni (1-step) ma non è una politica miope perché consente di ottenere una politica globalmente ottima.



Problematiche legate al calcolo di $V(s)$



Soluzione vicina alla ricerca esaustiva. Devo valutare per ogni stato tutte le possibili azioni (devo trovare il massimo).

Per tutte le possibili azioni devo calcolare la probabilità di transizione allo stato successivo e di ottenere una certa reward.

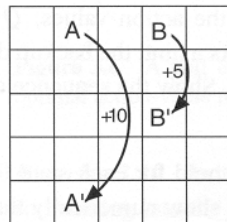
3 assunzioni:

- 1) Conoscenza della dinamica dell'ambiente: $P(s \rightarrow s' | a_j)$
- 2) Potenza di calcolo sufficiente
- 3) Proprietà Markoviane dell'ambiente (definizione di uno stato).

Soluzioni approssimate.



Esempio (possibile progetto)



(a)

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

(b)

Figure 3.5 Grid example: (a) exceptional reward dynamics; (b) state-value function for the equiprobable random policy.

$a = \pi(s)$ equiprobabile = $0.25 \forall a$

$R(s,a)$ per stati sui bordi e azioni che portano fuori dal bordo.

$R(s)$ per tutti gli altri stati ($= 0, +5, +10$)



Esempio di valore di uno stato

$t: s = s_t = C$

$t + 1: s = s' = s_{t+1}$

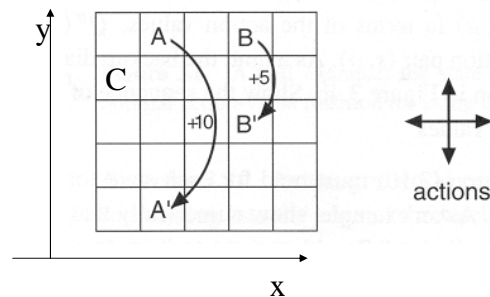
$C + Dy$ $p = 0.25$ $r = 0$

$C - Dy$ $p = 0.25$ $r = 0$

$C + Dx$ $p = 0.25$ $r = 0$

C $p = 0.25$ $r = -1$

$r_{t+1} = -0.25$



$t + 2, a = +Dy$

$\rightarrow C + Dy$ $p = 0.25$ $r = -1$

$\rightarrow C$ $p = 0.25$ $r = 0$

$\rightarrow A$ $p = 0.25$ $r = 0$

$\rightarrow C + Dy$ $p = 0.25$ $r = -1$

$r_{t+1} = -0.50$

$a = -Dy$

$\rightarrow C$ $r = 0$

$\rightarrow C - 2Dy$ $r = 0$

$\rightarrow C - Dy + Dx$ $r = 0$

$\rightarrow C - Dy$ $r = -1$

$r_{t+1} = -0.25$



Calcolo della funzione valore di uno stato



t: $s = s_t = C$

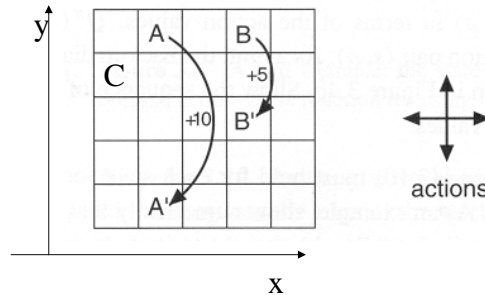
t + 1: $s = s' = s_{t+1}$

I $C + Dy$ $p = 0.25$ $r = 0$

II $C - Dy$ $p = 0.25$ $r = 0$

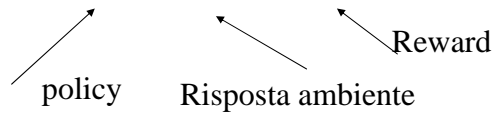
III $C + Dx$ $p = 0.25$ $r = 0$

IV C $p = 0.25$ $r = -1$



Per calcolare il valore di s : $V(s)$, devo analizzare il valore di ogni s' :

$\pi(s,a)$ -- $P_{s \rightarrow s' | a}$ -- $R_{s \rightarrow s' | a}$ -- $V^\pi(s')$ ← Value of s' .



A.A. 2006-2007

25/26

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



Sommario



Le equazioni di Bellman

Value function ottima

A.A. 2006-2007

26/26

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