

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. 2005-2006

1/22

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



Sommario



Le equazioni di Bellman

Value function ottima

A.A. 2005-2006

2/22

<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: $a = A(s)$

Policy di un agente: $\pi(s, a)$ è quanto dobbiamo definire.

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

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

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

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

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

A.A. 2005-2006

3/22

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



La value function



Nulla è detto sulla policy: dato uno stato, in quale nodo azione 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:

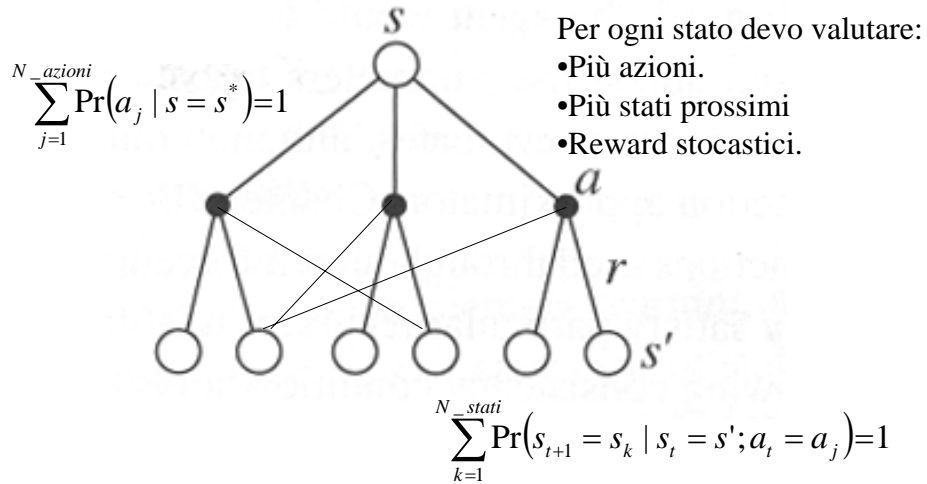
A.A. 2005-2006

4/22

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



Value function e modelli markoviani



A.A. 2005-2006

5/22

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



Calcolo ricorsivo della Value function



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

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

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

A.A. 2005-2006

6/22

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

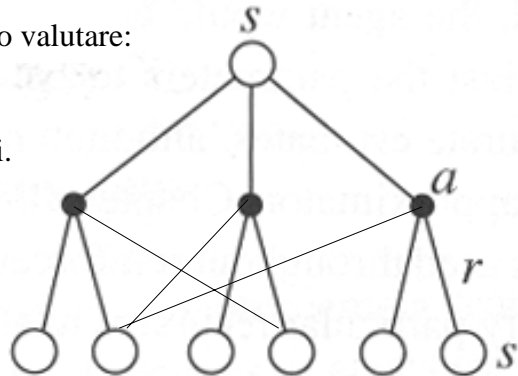


$V(s)$: primo termine

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

Per ogni stato devo valutare:

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



A.A. 2005-2006

7/22

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

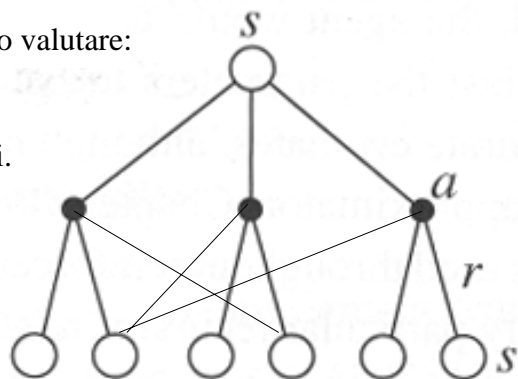


$V(s)$: secondo termine

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

Per ogni stato devo valutare:

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



A.A. 2005-2006

8/22

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



Calcolo ricorsivo della Value function



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

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

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

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

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

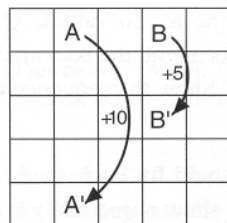
A.A. 2005-2006

9/22

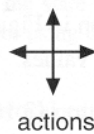
<http://homes.dsi.unimi.it/~borghese/>



Esempio (possibile progetto)



(a)



actions

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.A. 2005-2006

10/22

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



Esempio di valore di uno stato



t: $s = s_t = C$

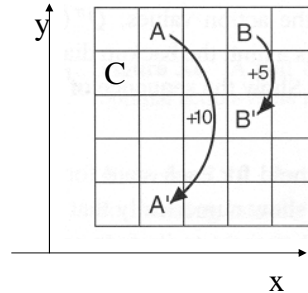
t + 1: $s = s' = s_{t+1}$

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



t + 2, a = +Dy

→ C + Dy p = 0.25 r = -1

→ C p = 0.25 r = 0

→ A p = 0.25 r = 0

→ C + Dy p = 0.25 r = -1

a = -Dy

→ C r = 0

→ C - 2Dy r = 0

→ C - Dy + Dx r = 0

→ C - Dy r = -1

A.A. 2005-2006

11/22

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



Calcolo della funzione valore di uno stato



t: $s = s_t = C$

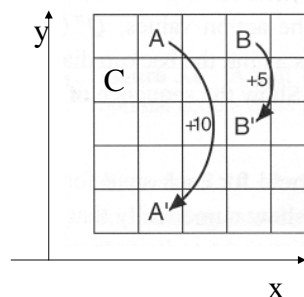
t + 1: $s = s' = s_{t+1}$

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' :

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

policy

Risposta ambiente

Reward

A.A. 2005-2006

12/22

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



Reinforcement Learning Problem



Given: Repeatedly...

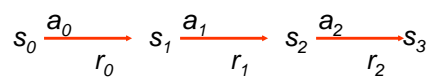
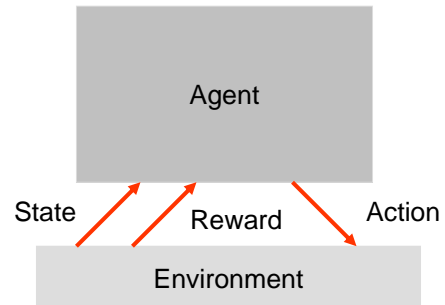
- Executed action
- Observed state
- Observed reward

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



How About Learning the Policy Directly?



1. $\pi^*: S \rightarrow A$
2. fill out table entries for π^* by collecting statistics on training pairs $\langle s, a^* \rangle$.
3. Where does a^* come from?

Problema di ottimizzazione della Value function, voglio determinare la policy che mi massimizza la Value function.



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:

$$V^*(s) \geq V^{\pi_k}(s) \quad \forall s.$$



Calcolo ricorsivo della Value function ottima



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

$V^*(s)$ 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 .

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

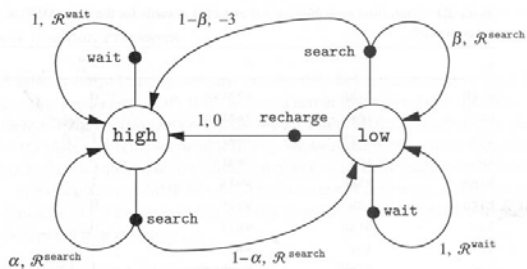
A.A. 2005-2006

17/22

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



Esempio di calcolo della funzione valore



$\alpha=0.6, \beta=0.1, \gamma=0.8, R_{\text{search}}=3, R_{\text{wait}}=1$

$\Pr(a, \text{state}=\text{high}) = [0.4, 0.6, 0]$

$\Pr(a, \text{state}=\text{low}) = [0.4, 0.5, 0.1]$

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

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

Sistema lineare di 2 equazioni nelle 2 incognite: V_h e V_l

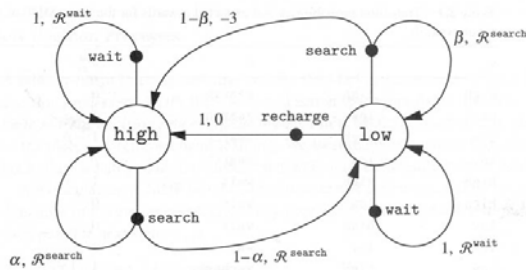
A.A. 2005-2006

18/22

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



Esempio di calcolo della policy ottima



V(high):

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

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

V(low):

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

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

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

A.A. 2005-2006

19/22

<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.

A.A. 2005-2006

20/22

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



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

Soluzioni approssimate.



Sommario



Le equazioni di Bellman

Value function ottima