

Policy Improvement

Alberto Borghese

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



A.A. 2013-2014

1/31

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



Sommario



Come migliorare la policy (Value iteration)

Esempi

A.A. 2013-2014

2/31

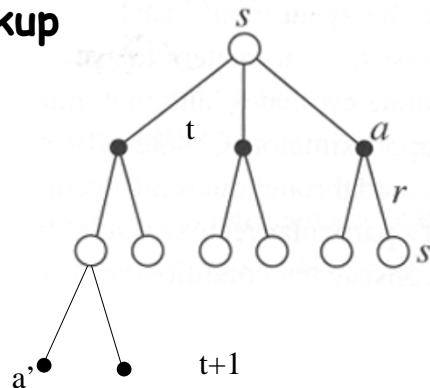
<http://\homes.dsi.unimi.it/~borghese/>



Tecnica full-backup

Back-up ↑

$\pi(s,a)$ fissata



Conosciamo $Q_k(s_t, a_t) \forall s_t$, anche per s'_{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 : $Q_{k+1}(s_t, a_t)$ congruente con:

$Q_k(s_{t+1}, a_{t+1})$ ed r_{t+1}

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

Da s' mi guardo indietro ed aggiorno $Q(s, a)$.

π fissata

A.A. 2013-2014

3/31

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



Calcolo iterativo della Value Function



Per ogni stato s , estratto a caso, analizziamo una singola transizione.

Equazione di Bellman per “**iterative policy evaluation**”:

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

Mi fido di $Q_{k+1}(s', a')$ (Backup)

$$\lim_{k \rightarrow \infty} \{Q_k(s, a)\} = Q^\pi(s, a)$$

A.A. 2013-2014

4/31

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

Iterative policy evaluation

Sommaruga_PolicyIteration_VisualBasic_Labirinto

A.A. 2013-2014 5/31 http://homes.dsi.unimi.it/~borghese/

Relazione soddisfatta da $Q^*(s, a)$

$$Q^*(s, a) = \underset{a_{t+1}}{\text{Max}} [E_{\pi} \{ R_t \mid s_t = s, a_t = a \}] =$$

$$\underset{a_{t+1}}{\text{Max}} \left[E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\} \right] =$$

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

$$\underset{a_{t+1}}{\text{Max}} [r_{t+1} + \gamma Q^*(s_{t+1}, a_{t+1}) \mid s_t = s, a_t = a] \Rightarrow$$

$$Q^*(s, a) = \underset{a'}{\text{Max}} \left[P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + Q^*(s', a')] \right]$$

Bellman's Equation For optimal policy

A.A. 2013-2014 6/31 http://homes.dsi.unimi.it/~borghese/



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:

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

$$\sum_{s'} P_{s_m \rightarrow s'}^{a_{new}} [R_{s_m \rightarrow s'}^{a_{new}} + \gamma V^\pi(s')] \quad V(s) = \text{value function sullo stato}$$

?

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

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” ad un passo (reward del passo + reward successivo).



Enunciato del teorema del miglioramento della policy

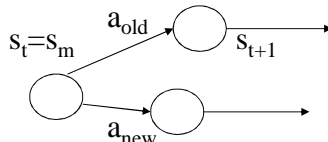


$$Q^\pi(s, a) = \sum_k P_{s \rightarrow s_k | a} [R_{s \rightarrow s_k | a} + \gamma V^\pi(s_k)]$$

Ipotesi: π and π' deterministi 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: $V^{\pi'}(s) \geq V^\pi(s) \forall s$.
 $Q^{\pi'}(s, a_{new}) \geq Q^\pi(s, a_{old})$



A.A. 2013-2014

9/31

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



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. 2013-2014

10/31

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



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

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

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

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

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

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

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



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

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

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

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

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



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:

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

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



Calcolo ricorsivo della Value function ottima: confronti



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

$Q^*(s,a)$ 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 .

$$Q^*(s,a) = \max_{a'} \sum_{s'} P_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma Q^*(s', a')]$$

Politica greedy: scelgo l'azione ottimale.

Ha senso per il robot raccogli-lattine?



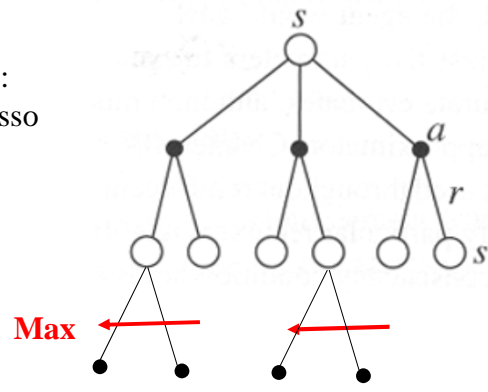
Q*(s, a) - Osservazioni



$$Q^*(s, a) = \max_{a'} \sum_{s'} P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma Q^*(s', a')]$$

Per ogni stato devo valutare:
 • L'azione migliore ad un passo

Come valuto?
 • analizzando reward a lungo termine



$t+1$



Policy iteration



Iterazione tra:

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

$$\pi_0 \rightarrow V^{\pi_0} \rightarrow \pi_1 \rightarrow V^{\pi_1} \rightarrow \pi_2 \rightarrow V^{\pi_2} \rightarrow \dots$$

$$\rightarrow \pi^* \rightarrow V^*$$

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



Algoritmo - I



Inizialization

$Q(s,a) = 0$;
 $\pi(s,a) = \text{random}$ (e.g. equiprobabile);

Repeat
 point 2.
 point 3.
until policy_stable



Algoritmo - II



Policy evaluation – versione per trial

Repeat
 Th = 0; // small value;
 for s = 1:N
 for a = 1:M
 $Q_temp = \sum_{s'} \Pr_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma \sum_{a'} \pi(s', a') Q(s', a')]$
 $\Delta Q = |Q(s,a) - Q_temp|$
 $Q(s,a) = Q_temp$;
 th = max(th, ΔQ)
 end;
 end;
until th < th_max;



Algoritmo - III



Policy improvement

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



Algoritmo - II



Policy evaluation – versione per epoch

```
Repeat
    Th = 0; // small value;
    for s = 1:N
        for a = 1:M
             $Q\_temp(s,a) = \sum_{s'} \Pr_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma \sum_{a'} \Pr_{a'|s'} Q(s', a')]$ 
             $\Delta Q = |Q(s,a) - Q\_temp(s,a)|$ 
            th = max(th,  $\Delta Q$ )
        end;
    end;
    for s = 1:N, for a=1:m
         $Q(s,a) = Q\_temp(s,a)$ ;
    end; end;
until th < th_max;
```



Max or soft max



Policy improvement

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

Max con policy ϵ -greedy, soft-max, ...



Iterative policy evaluation - problema



$$V_{k+1}(s) = \left[\sum_{a_j} \pi(a_j, s) \right] \sum_{s'} P_{s \rightarrow s'|a_j} [R_{s \rightarrow s'|a_j} + \gamma V_k(s')]$$

Converge al limite a $V^\pi(s)$. Come facciamo a troncare?



Value iteration



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

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

$$Q_{k+1}(s, a) = \max_{a'} \sum_{s'} P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma Q_k(s', a')]$$

$\forall s, a$

A.A. 2013-2014

23/31

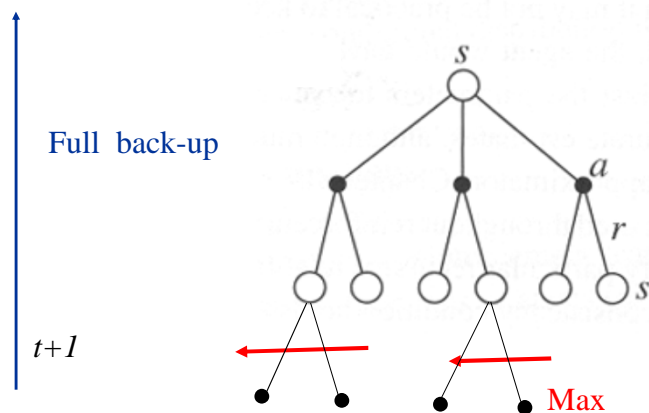
<http://borghese.di.unimi.it/>



Visualizzazione grafica



$$V_{k+1}(s) = \max_a \sum_{s'} P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V_k(s')]$$



A.A. 2013-2014

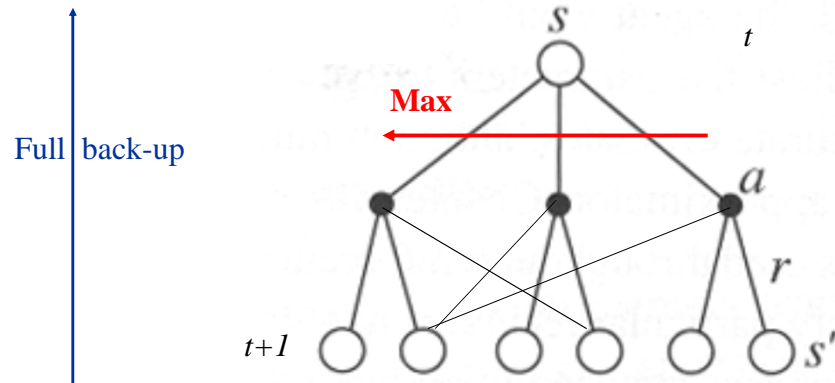
24/31

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



Visualizzazione grafica

$$V_{k+1}(s) = \max_a \sum_{s'} P_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma V_k(s')]$$



A.A. 2013-2014

25/31

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



Sommario

Come migliorare la policy (Value iteration)

Esempi

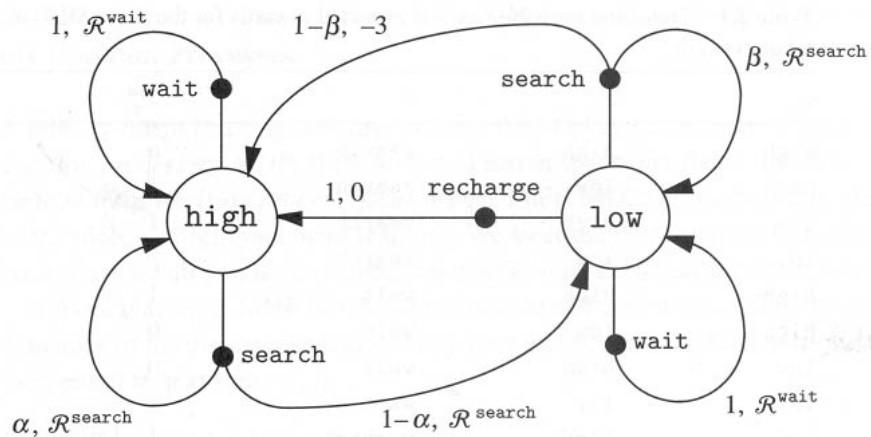
A.A. 2013-2014

26/31

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



Robot cerca-lattine



A.A. 2013-2014

27/31

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



Esempio: robot - Policy deterministica



$$Q(h, \text{search}) = \Pr(h \rightarrow l, \text{search}) \times [R(h \rightarrow h, \text{search}) + \gamma \times Q(h, \text{search})]$$

$$\Pr(h \rightarrow h, \text{search}) \times [R(h \rightarrow l, \text{search}) + \gamma \times Q(l, \text{wait})]$$

$$Q(h, \text{search}) = 0.4 \times [3 + 0.8 \times Q(h, \text{search})] + 0.6 \times [3 + 0.8 \times Q(l, \text{wait})]$$

$$Q(l, \text{wait}) = \Pr(l \rightarrow l, \text{wait}) \times [R(l \rightarrow l, \text{wait}) + 0.8 \times Q(l, \text{wait})]$$

$$Q(l, \text{wait}) = 1 \times [1 + 0.8 \times Q(l, \text{wait})]$$

Policy iniziale deterministica:

STATO: Q(h,search) →

$$Q(h, s) \cong 4,4 + 0.7 \times Q(l, w) \cong 7.95$$

STATO: Q(l, wait) →

$$Q(l, \text{wait}) = 5$$



Posso migliorare la policy?

A.A. 2013-2014

28/31

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



Esempio: robot - miglioramento policy



Miglioro la policy, modificando l'azione associata a $s = \text{low}$:

STATO: high

$a = \text{search} \rightarrow Q(\text{h}, \text{search}) \cong 4.4 + 0.7 Q(\text{l}, \text{recharge}) \neq 7.95$

STATO: low

$a = \text{recharge} \rightarrow Q(\text{l}, \text{recharge}) = 0 + 0.8 Q(\text{h}, \text{search}) = ???$

Ho stimato correttamente $Q(\text{h}, \text{search})$? No

Applico iterative policy evaluation



STATO: VI

$a = \text{recharge} \rightarrow Q_1(\text{l}, \text{r}) = 0.8 Q_1(\text{h}, \text{s}) = 0.8 \times 7.95 = 6.36$

STATO: high

$a = \text{search} \rightarrow Q_2(\text{h}, \text{s}) \cong 4.4 + 0.7 Q_1(\text{l}, \text{r}) \cong 4.4 + 0.7 \times 6.36 = 8.85$

Ho stimato correttamente $Q(\text{s}, \text{a})$? No. Devo iterare la policy evaluation.



Esempio: robot - IV



Asintoticamente calcolo il valore vero delle coppie stato-azione:

STATO: high

$a = \text{search} \rightarrow Q(\text{h}, \text{s}) \cong Q_2(\text{h}, \text{s}) \cong 4.4 + 0.7 Q_1(\text{l}, \text{r}) = 4.4 + 0.7 \times 6.36 = 8.85$

STATO: low

$a = \text{recharge} \rightarrow Q(\text{l}, \text{r}) = 0.8 Q(\text{h}, \text{s}) \rightarrow 7.1$

Potrei ottenere gli stessi valori ottenuti asintoticamente, risolvendo il sistema lineare:

$$Q(\text{h}, \text{s}) = 4.4 + 0.7 Q(\text{l}, \text{r}) =$$

$$Q(\text{l}, \text{r}) = 0.8 Q(\text{h}, \text{s}) =$$

Ho terminato?





Sommario



Come migliorare la policy (Value iteration)

Esempi