

Sistemi Intelligenti Reinforcement Learning: Iterative policy evaluation

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. 2011-2012

1/40

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



Sommario



Determinazione ricorsiva della Value function

Determinazione della policy ottima.

A.A. 2011-2012

2/40

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

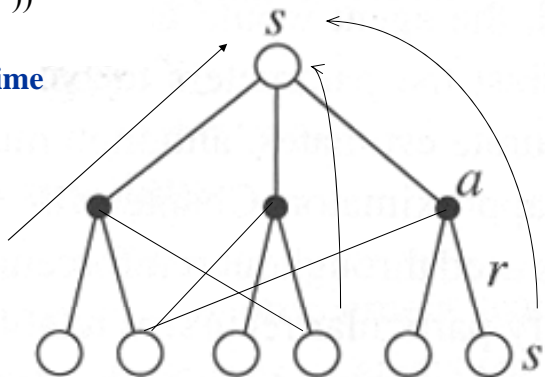


Osservazioni



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

Backwards in time



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

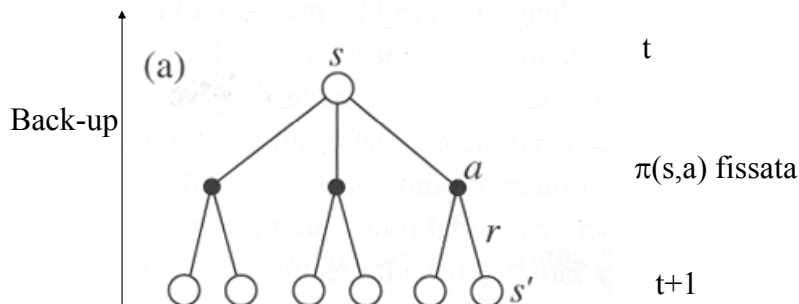
A.A. 2011-2012

3/40

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



Tecnica full-backup



Conosciamo $V_k(s(t)) \forall s$, anche per s' quindi

Analizziamo la transizione da $s(t) \rightarrow \{s'(t+1)\}$

Calcoliamo un nuovo valore di s : $V_{k+1}(s(t))$ congruente con $V_k(s(t))$ ed r_{t+1}

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

π fissata

Da s' mi guardo indietro ed aggiorno $V(s)$.

A.A. 2011-2012

4/40

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



Calcolo ricorsivo della Value function



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

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

Legame?

Policy Next-state

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

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

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



Calcolo iterativo della Value Function



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

Equazione di Bellman per “iterative policy evaluation”:

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

Mi fido di $V_k(s')$ (Backup)

$$\lim_{k \rightarrow \infty} \{V_k(s)\} = V^\pi(s)$$



Iterative policy evaluation



Evoluzione del sistema da $s(t=0)$ a $\{s'(t=T)\}$ utilizzando la policy $\pi(s,a)$, prefissata.

Quanto valgono gli stati?

Parto da $V(s(t=0))_{k=0}$ arbitraria, otterrò una value function per ogni stato che sarà funzione di $V(s(t=0))$.

Devo migliorare, come?

Utilizziamo l'informazione sul **passato**.

$$\{V\}^0, \{V\}^1, \{V\}^2, \{V\}^3, \{V\}^4, \{V\}^5, \dots \{V\}^\infty$$
$$\lim_{k \rightarrow \infty} \{V_k(s)\} = V^\pi(s)$$



Fondamenti del metodo



- Supponiamo di essere all'istante t . In questo istante t , si può passare ad 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.



Algoritmo per "iterative policy evaluation", versione batch



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

Definiamo una soglia di convergenza τ

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

Repeat

```

{
  Δ = 0;
  for s = 1 : N
    // ∀ s, ≠ TS
    {
      W(s) =  $\sum_{a_j} \pi(s, a_j) \sum_{s'} P_{s \rightarrow s'}^{a_j} [R_{s \rightarrow s'}^{a_j} + \gamma V(s')]$  // W(s) è V_{k+1}(s)
      Δ = max(Δ, |V(s) - W(s)|)
    }
  }
  for s=1:N
    V(s) = W(s);
} Until (Δ < τ);

```

A.A. 2011-2012

9/40

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



Interpretazione dell'update (batch o trial)



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

Al termine dell'aggiornamento dei $V(s)$ per tutti gli stati, $V(s) = V_{\text{new}}(s)$. **Aggiornamento batch.**

Utilizzerò in parte già il nuovo valore di $V(s)$ all'interno dell'equazione di aggiornamento. **Aggiornamento per trial.**

Entrambe le modalità di aggiornamento convergono.

A.A. 2011-2012

10/40

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



Algoritmo per "iterative policy evaluation", versione per trial



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

Definiamo una soglia **relativa** di convergenza τ

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

Repeat

{ $\Delta = 0;$

for $s = 1 : N$ // $\forall s, \neq TS$

{ Value = $V(s);$

$$V_{k+1}(s) = \sum_{a_j} \pi(s, a_j) \sum_{s'} P_{s \rightarrow s'}^{a_j} [R_{s \rightarrow s'}^{a_j} + \gamma V(s')]$$

Forward
pass

$$\Delta = \max(\Delta, (| \text{Value} - V_{k+1}(s) |) / | \text{Value} |)$$

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



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



3 assunzioni:

- 1) Conoscenza della dinamica dell'ambiente: $P(s \rightarrow s' | a_j)$
- 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).



Sommario



Determinazione ricorsiva della Value function

Determinazione della policy ottima.



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.



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.

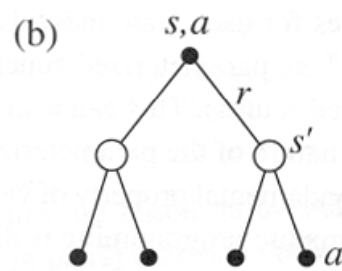
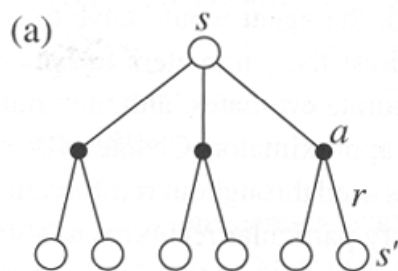
3 assunzioni:

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

Soluzioni approssimate.



Policy



La policy deve essere ancora determinata. Come fa l'agente a determinare la policy ottimale?

Archi multipli fuoriuscenti da un'azione sono associati alla probabilità di scegliere quel cammino (ambiente stocastico).

Archi multipli fuoriuscenti da uno stato, sono associati alla policy.



Relazione soddisfatta da $V^*(s)$



$$V^*(s) = \underset{a}{\text{Max}} [E_{\pi} \{R_t | s_t = s, a_t\}] =$$

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

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

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

$$V^*(s) = \underset{a}{\text{Max}} \{ P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V^*(s')] \}$$

Bellman's
Equation
For optimal
policy

A.A. 2011-2012

17/40

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

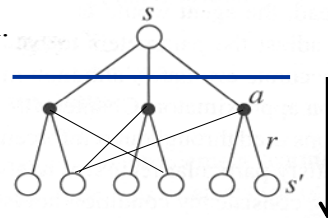


La Value function sulle azioni



La value function può riguardare le azioni.

Action-Value function



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

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

Dove abbiamo già trovato una Value function associata alle azioni?

A.A. 2011-2012

18/40

<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 ? $a' \neq \pi(s_m)$.

Cosa viene influenzato?

Scelgo a' in s , 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 , $\{s_{t+1}\}$, ed il reward a lungo termine:

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

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

?

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

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

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

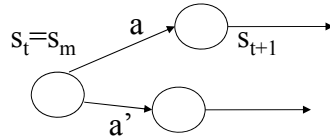


$$Q^{\pi}(s, a) = \sum_k P_{s \rightarrow s_k | a} [R_{s \rightarrow s_k | a} + \gamma V^{\pi}(s_k)]$$

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

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

Tesi: π' è meglio di π . Cioè: $V^{\pi'}(s) \geq V^{\pi}(s) \forall s$.



A.A. 2011-2012

21/40

<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. 2011-2012

22/40

<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

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

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

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

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

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

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



Ottimizzazione 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: confronti



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

$V^*(s)$ 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 .

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

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



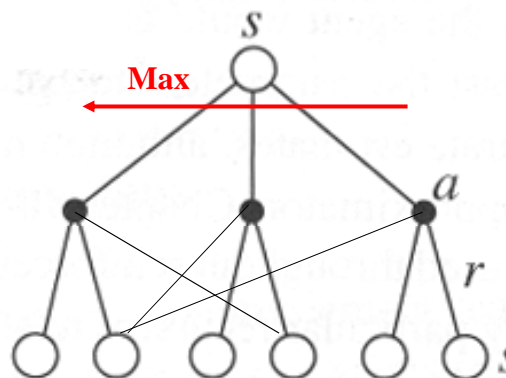
$V^*(s)$ - Osservazioni



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

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

Come valuto?
• analizzando reward a lungo termine





Utilizzo di $V(s)$ per determinare la policy ottima



Esplorazione dell'effetto a lungo termine di tutte le azioni possibili.

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.

E' reso possibile per la conoscenza dell'ambiente (stocastico).



Politica ottima



Miglioramento della politica per tutti gli stati.

$$\begin{aligned} \pi'(s_k) &= \arg \max_a Q^\pi(s, a) && \text{greedy o } \varepsilon\text{-greedy} \\ &= \arg \max_a E\{r_{t+1} + \gamma V^\pi(s') \mid s_t = s, a_t = a\} \\ \forall s & \\ &= \arg \max_a \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma V^\pi(s')] \end{aligned}$$

Policy improvement

Si può estendere al caso di comportamento stocastico dell'agente nel qual caso: $\pi(s,a)$ è una probabilità. Questo sarà utile quando l'ambiente non sarà più considerato noto.



Policy ottima

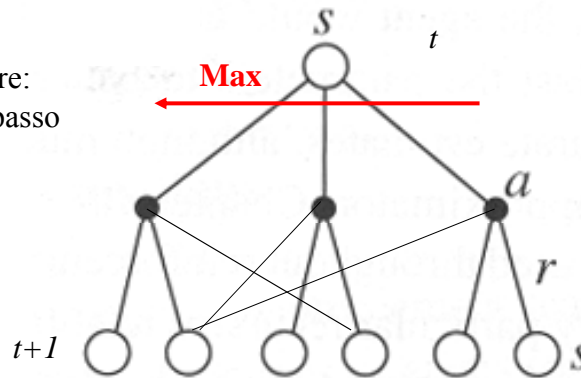


$$V^*(s) = \max_{a_j} \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma V^*(s')]]$$

$$\pi^*(s) = \arg \max_a Q^\pi(s, a) \quad \forall s$$

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

Come valuto?
 • analizzando reward a lungo termine



A.A. 2011-2012

29/40

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



Policy iteration



Iterazione tra:

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

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

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

Converge velocemente ad una buona politica

A.A. 2011-2012

30/40

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



Algoritmo (progetto per esame)



Repeat until
policy-stable

1. Initialization
 $V(s) \in \mathfrak{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$
2. Policy Evaluation
Repeat
 $\Delta \leftarrow 0$
For each $s \in \mathcal{S}$:
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{s'} \mathcal{P}_{ss'}^{\pi(s)} [\mathcal{R}_{ss'}^{\pi(s)} + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$ (a small positive number)
3. Policy Improvement
 $policy\text{-}stable \leftarrow true$
For each $s \in \mathcal{S}$:
 $b \leftarrow \pi(s)$
 $\pi(s) \leftarrow \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$
If $b \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
If $policy\text{-}stable$, then stop; else go to 2



Politica ottima



Miglioramento della politica per tutti gli stati.

$$\begin{aligned}
 \pi^*(s_k) &= \arg \max_a Q^\pi(s, a) \\
 &= \arg \max_a Q^\pi(s, a) \\
 \forall s &= \arg \max_a E\{r_{t+1} + \gamma \mathcal{V}^\pi(s') \mid s_t = s, a_t = a\} \\
 &= \arg \max_a \sum_{s'} P_{s \rightarrow s'}^a [\mathcal{R}_{s \rightarrow s'}^a + \gamma \mathcal{V}^\pi(s')]
 \end{aligned}$$

Si può estendere al caso di comportamento stocastico dell'agente
nel qual caso: $\pi(s,a)$ è una probabilità.



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



Iterative policy evaluation - problema



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

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



Value iteration

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

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

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

$\forall s$

A.A. 2011-2012

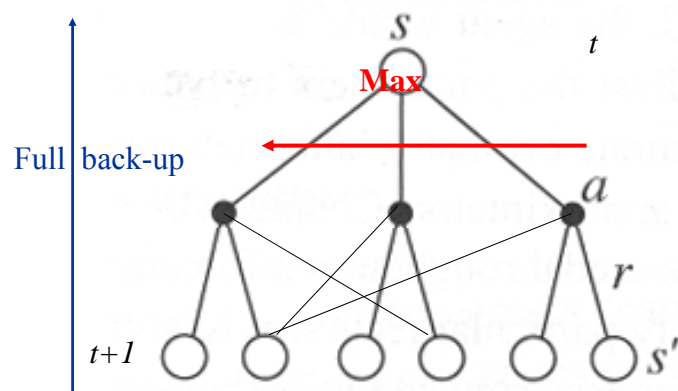
35/40

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



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. 2011-2012

36/40

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



Confronto con l'equazione di Bellman



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

$V^*(s)$ 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 .

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

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

A.A. 2011-2012

37/40

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



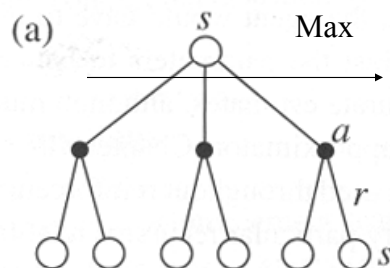
Confronto con l'equazione di backup



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

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

Nel caso della Value Iteration viene considerata solamente l'azione che fornisce il valore massimo.



A.A. 2011-2012

38/40



Algoritmo di value iteration



$V(s) = 0 \forall s$ compreso TS

Repeat

```
{
  Δ = 0;
  for s = 1:NS // Per ogni stato eccetto il TS
  {
    V_buffer = V(s);
    V(s) = max_a { ∑_{s'} P_{s→s'|a} [R_{s→s'|a} + γV(s')] }
    Δ = max(Δ, |V(s) - V_buffer|);
  }
} until (Δ < Th);
```

Output a deterministic policy such that:

$$\pi(s) = \arg \max_a \sum_{s'} P_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma V(s')]$$



Sommario



Determinazione ricorsiva della Value function

Determinazione della policy ottima.