

Sistemi Intelligenti Reinforcement Learning: Value Iteration and Temporal Differences

Alberto Borghese

Università degli Studi di Milano
Laboratorio di Sistemi Intelligenti Applicati (AIS-La)
Dipartimento di Informatica
borgnese@di.unimi.it



A.A. 2019-2020

1/37

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



Sommario



Value iteration

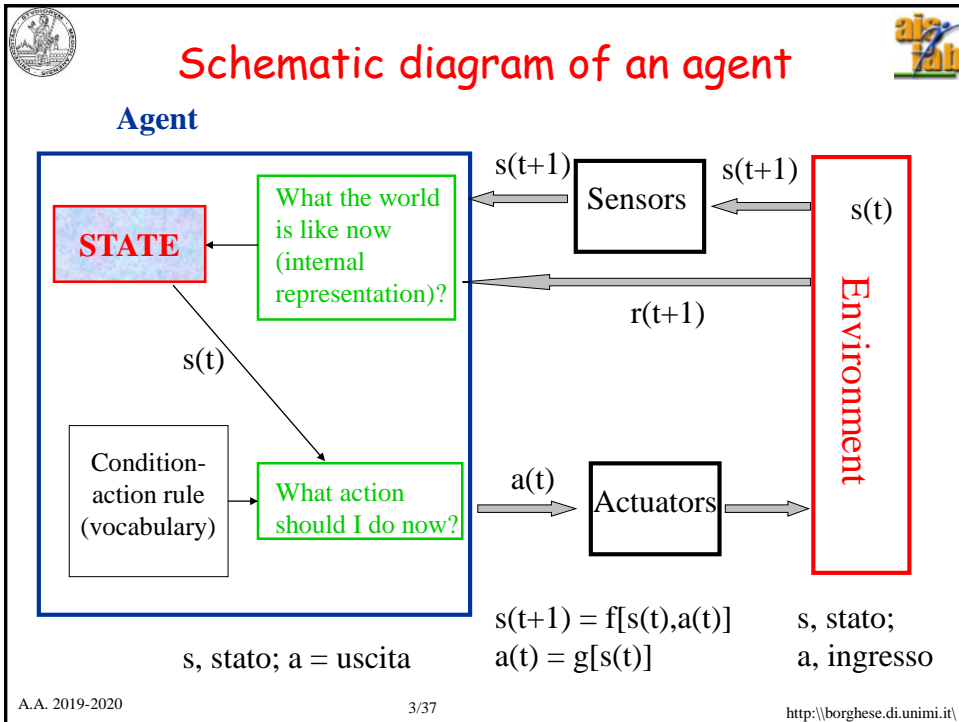
Esempi

Learning with temporal differences

A.A. 2019-2020

2/37

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



Calcolo ricorsivo della Value function ottima: confronti

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

$Q^*(s, a)$ 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 .

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

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

A.A. 2019-2020 4/37 <http://borghese.di.unimi.it/>



$V^*(s)$ - Osservazioni



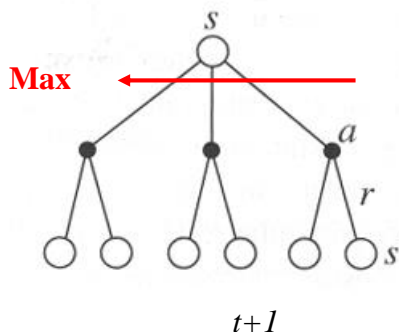
$$V^*(s) = \max_a \sum_{s'} P_{s \rightarrow s'|a} [R_{s \rightarrow s'|a} + \gamma V^*(s')]$$

Per ogni stato devo valutare:

- L'azione migliore ad un passo

Come valuto?

- analizzando reward a lungo termine



Policy iteration



Iterazione tra:

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

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

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



Algoritmo



Inizialization

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

Repeat
 point 2.
 point 3.
until policy_stable



Algoritmo - point2



Policy evaluation – versione per trial

Repeat
 th = 0; // small value;
 for s = 1:N

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

 $\Delta V = |V(s) - V_temp|$
 $V(s) = V_temp$;
 th = max(th, ΔV)
 end;
until th < th_max;



Algoritmo - point3



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 V(s')] \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

$$V_temp(s,a) = \sum_{a_j} \pi(s, a_j) \sum_{s'} \Pr_{s \rightarrow s'|a_j} [R_{s \rightarrow s'|a_j} + \gamma V(s')]$$

$$\Delta V = |V(s) - V_temp(s)|$$

$$th = \max(th, \Delta V)$$

end;

end;

for s = 1:N

$$V(s) = V_temp(s);$$

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);

```

$$a_{\text{new}} = \underset{a}{\operatorname{arg\,max}} \left\{ \sum_a \pi(s, a) \sum_{s'} \Pr_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V(s')] \right\}$$

```

  if (a_new ≠ a_old)
    policy_stable = false;
  end;

```

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



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



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

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



Value iteration

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

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

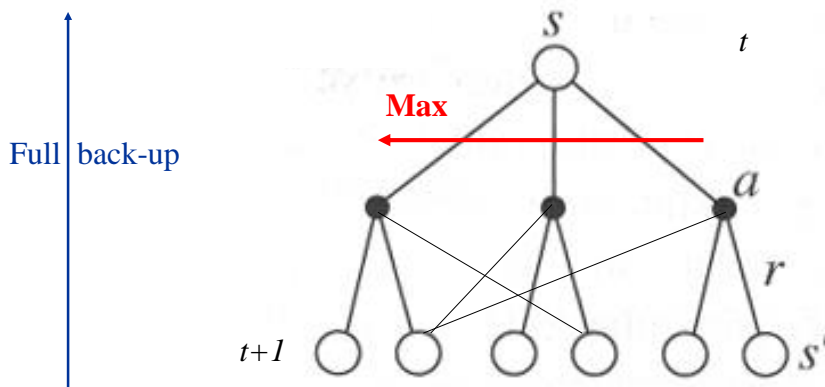
$$V_{k+1}(s) = \max_a \sum_{s'} \pi(a, s) P_{s \rightarrow s' | a} [R_{s \rightarrow s' | a} + \gamma V_k(s')]$$

$\forall s$



Visualizzazione grafica

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





Sommario

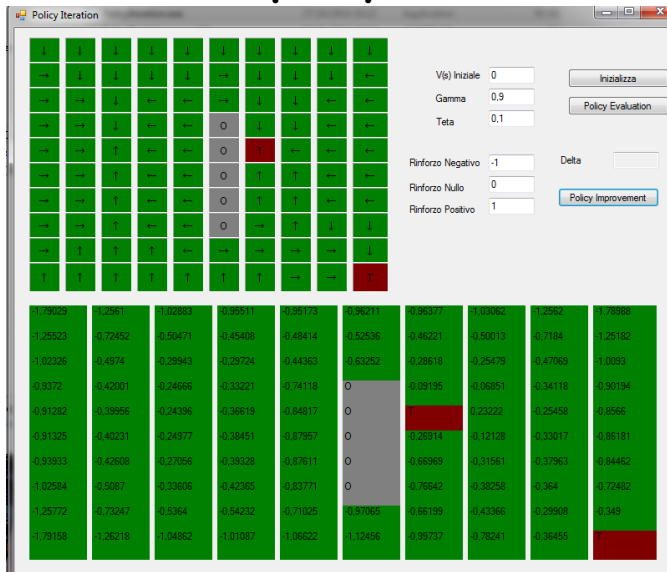
Value iteration

Esempi

Learning with temporal differences

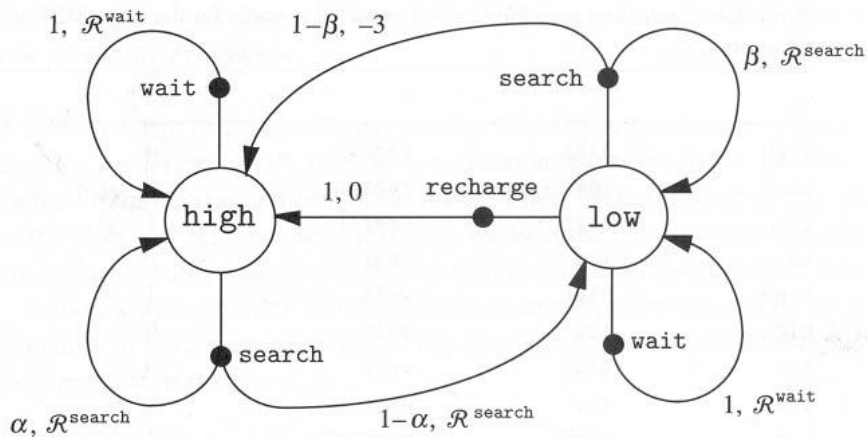


Iterative policy evaluation





Robot cerca-lattine



A.A. 2019-2020

17/37

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



Esempio: robot - Policy deterministica



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

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

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

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

Policy iniziale deterministica:

STATO: $Q(h,search) \rightarrow$
 $Q(h,s) \cong 4,4 + 0.7 \times Q(l,w) \cong 7.95$
STATO: $Q(l, wait) \rightarrow$
 $Q(l,wait) = 5$



Posso migliorare la policy?

A.A. 2019-2020

18/37

<http://borgese.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 un passo di iterative policy evaluation, in modalità trial.



STATO: VI

$$a = \text{recharge} \rightarrow Q(\text{l}, \text{r}) = 0,8 Q(\text{h}, \text{s}) = 0,8 \times 7,95 = 6,36$$

STATO: high

$$a = \text{search} \rightarrow Q(\text{h}, \text{s}) \cong 4,4 + 0,7 Q(\text{l}, \text{r}) \cong 4,4 + 0,7 \times 6,36 = 8,85$$

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

$$a = \text{search} \rightarrow Q(\text{h}, \text{s}) \cong 4,4 + 0,7 Q_1(\text{l}, \text{r}) = 4,4 + 0,7 \times 6,36 = 8,85$$

STATO: low

$$a = \text{recharge} \rightarrow Q(\text{l}, \text{r}) = 0,8 Q(\text{h}, \text{s}) = 0,8 \times 8,85 = 7,08$$

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

$$Q(\text{h}, \text{s}) = 4,4 + 0,7 Q(\text{l}, \text{r}) = 10$$

$$Q(\text{l}, \text{r}) = 0,8 Q(\text{h}, \text{s}) = 8$$

A questi valori si arriva solo asintoticamente





Sommario



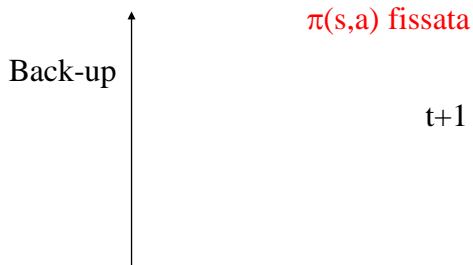
Value iteration

Esempi

Learning with temporal differences



Tecnica full-backup



Conosciamo $V_k(s_t) \forall s_t$, anche per s'_{t+1} quindi:

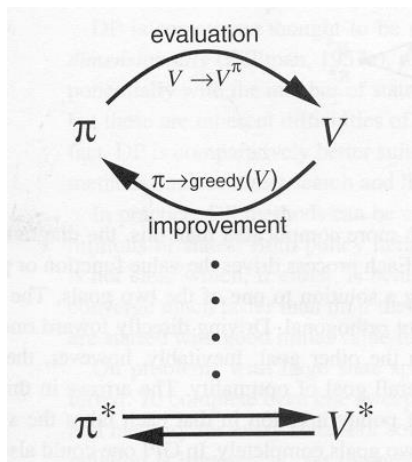
Analizziamo la transizione da $s_t, a_t \rightarrow (s'_{t+1})$

Calcoliamo un nuovo valore di V per s : $V_{k+1}(s_t)$ congruente con:

$$V_k(s_{t+1}) \text{ ed } r_{t+1}$$

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

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



Generalized Policy iteration

Schema di Apprendimento

Policy iteration
Value iteration

Competizione e cooperazione -> V corretta e policy ottimale.



World RL competition

Started in NIP2006.

It became very popular and started a workshop on its own.
Last official challenge at COLT was at 2014. Still on-going
research and specific challenges.

Visit: <http://www.rl-competition.org/>

Project Malmo competition.

<https://www.microsoft.com/en-us/research/project/project-malmo/> (2019)



How About Learning the Value Function?



Facciamo imparare all'agente la value function, per una certa politica: V^π :

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

È una funzione dello stato.

Una volta imparata la value function, V^* , l'agente seleziona la policy ottima passo per passo, "one step lookahead":

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

Full backup, for all states



Value iteration



Facciamo imparare all'agente la value function, per una certa politica: V^π , analizzando quello che succede in uno step temporale:

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

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

L'apprendimento della policy si può inglobare nella value iteration:

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

vs



Problema legato alla conoscenza della risposta dell'ambiente



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

Full backup, single state, s , all future states s'

Fino a questo punto, è noto un modello dell'ambiente:

- $R(\cdot)$
- $P(\cdot)$

Environment modeling -> Value function computation -> Policy optimization.



Osservazioni



Iterazione tra:

- Calcolo della Value function

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

- Miglioramento della policy

$$= \arg \max_{a_j} \sum_{s'} P_{s \rightarrow s' | a_j} [R_{s \rightarrow s' | a_j} + \gamma \mathcal{W}^\pi(s')]$$

Non sono noti



Background su Temporal Difference (TD) Learning



Al tempo t abbiamo a disposizione:

$$r_{t+1} = r, \quad R_{s \rightarrow s' | a_j}$$

$$s_{t+1} = s, \quad P_{s \rightarrow s' | a_j}$$

Reward certo

Transizione certa

vengono misurati dall'ambiente

Come si possono utilizzare per apprendere?



Confronto con il setting associativo



$$Q_{k+1} = Q_k - \frac{Q_k}{N_{k+1}} + \frac{r_{k+1}}{N_{k+1}} = Q_k + \alpha [r_{k+1} - Q_k]$$

Occupazione di memoria minima: Solo Q_k e k .

NB k è il numero di volte in cui è stata scelta a_j .

Questa forma è la base del RL. La sua forma generale è:

$$\text{NewEstimate} = \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

$$\text{NewEstimate} = \text{OldEstimate} + \text{StepSize} * \text{Error}.$$

$$\text{StepSize} = \alpha = 1/k \quad a = \text{cost}$$

$$\text{Rewards weight } w = 1 \quad \text{Weight of } i\text{-th reward at time } k: w = (1-\alpha)^{k-i}$$

Qual è la differenza introdotta dall'approccio DP?



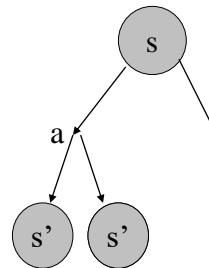
Un possibile aggiornamento

In iterative policy evaluation ottengo questo aggiornamento:

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

Ad ogni istante di tempo di ogni trial aggiorno la Value function:

$$V_{k+1}(s) = [r' + \gamma V_k(s)]$$



Qual'è il problema?



Un possibile aggiornamento di $Q(s, a)$

$$Q_{k+1} = Q_k - \frac{Q_k}{N_{k+1}} + \frac{r_{k+1}}{N_{k+1}} = Q_k + \alpha [r_{k+1} - Q_k] = Q_k + \Delta Q_k$$

Quanto vale α ?

$$V_k(s) = V_k(s) + \Delta V_k(s)$$

Come calcolo $\Delta V_k(s)$?



TD(0) update

Ad ogni istante di tempo di ogni trial aggiorno la Value function:

$$V_{k+1}(s_t) = V_k(s_t) + \alpha [r_{t+1} + \gamma V_k(s_{t+1}) - V_k(s_t)]$$

Da confrontare con la iterative policy evaluation:

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

Sample
backup

E con il valore di uno stato sotto la policy $\pi(s,a)$:

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

Quanto vale α ?



Confronto con il setting associativo

$$Q_{k+1} = Q_k - \frac{Q_k}{N_{k+1}} + \frac{r_{k+1}}{N_{k+1}} = Q_k + \alpha [r_{k+1} - Q_k]$$

Occupazione di memoria minima: Solo Q_k e k .
NB k è il numero di volte in cui è stata scelta a_j .

Questa forma è la base del RL. La sua forma generale è:

$$\begin{aligned} \text{NewEstimate} &= \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}] \\ \text{NewEstimate} &= \text{OldEstimate} + \text{StepSize} * \text{Error}. \end{aligned}$$

$$\text{StepSize} = \alpha = 1/k \quad a = \text{cost}$$



Setting α value



$\alpha(s_t, a_t, s_{t+1}) = 1/k(s_t, a_t, s_{t+1})$, where k represents the number of occurrences of s_t, a_t, s_{t+1} . With this setting the estimated Q tends to the expected value of $Q(s,a)$.

Per semplicità si assume solitamente $\alpha < 1$ costante. In questo caso, $Q(s,a)$ assume il valore di una media pesata dei reward a lungo termine collezionati da (s,a) , con peso: $(1-\alpha)^k$: *exponential recency-weighted average*.



Esempio: valutazione della policy TD



Stato	Tempo percorrenza stimato del segmento	Tempo percorrenza attuale del segmento	Tempo totale previsto in precedenza $V^k(s)$	Tempo totale previsto aggiornato $V^{k+1}(s)$	Increase or decrease $V(s)$	
S0	Esco dall'ufficio	0	0	30	$30 + (10 + 25 - 30) = 35$	>
S1	Salgo in auto	5	10	25	$25 + (15 + 10 - 25) = 25$	=
S2	Esco dall'autostrada	15	15	10	$10 + (10 + 5 - 10) = 15$	>
S3	Esco su strada secondaria	5	10	5	$5 + (2 + 3 - 5) = 5$	=
S4	Entro Strada di casa	2	2	3	$3 + (3 + 0 - 3) = 3$	=
S5	Parcheggio	3	3	0	0	

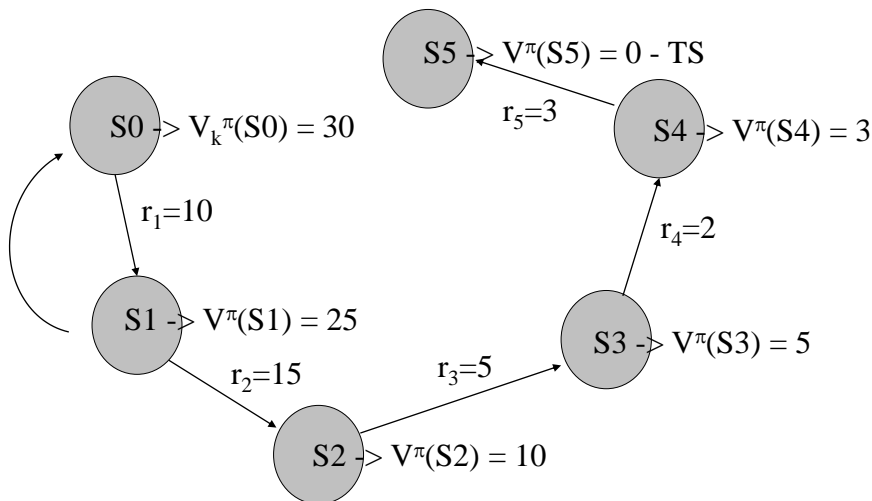
$V(s)$ è l'expected "Time-to-Go" - $\gamma = 1$ $\alpha = 1$



Learning $V_{\pi}(s)$



$$S0 \rightarrow V_{k+1}^{\pi}(S0) = 30 + (10+25-30) = 35$$



Come i diversi reward istantanei modificano $V^{\pi}(s)$?

A.A. 2019-2020

37/37

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



Alcuni passi di apprendimento di TD(0)



Alcuni passi di iterazione per TD(0)

$$V(0) = V(0) + \alpha (r_1 + \gamma V(1) - V(0)) = 30 + \alpha (5 + 35 - 30) = 30 + \alpha * \Delta$$

Stima iniziale del tempo di percorrenza totale: 30m

Tempo di percorrenza fino all'auto: 5m

Stima del tempo di percorrenza dal parcheggio: 35m

$$V(1) = V(1) + \alpha (r_1 + \gamma V(2) - V(1)) = 35 + \alpha (20 + 15 - 35) = 35 + \alpha * \Delta$$

Stima iniziale del tempo di percorrenza dal parcheggio: 35m

Tempo di percorrenza fino ad uscita autostrada: 20m

Tempo di percorrenza fino ad uscita autostrada: 20m

Stima del tempo di percorrenza dall'uscita autostrada: 15m

A.A. 2019-2020

38/37

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



Ruolo di α



$$V(1) = V(1) + \alpha (r_1 + \gamma V(2) - V(1)) = 30 + \alpha (10 + 25 - 30) = 30 - \alpha * 5$$

Stima iniziale del tempo di percorrenza dal parcheggio: 30m

Tempo per raggiungere l'auto: 10m

Stima del tempo di percorrenza dall'uscita Dal parcheggio: 25m

$\alpha < 1$.

If $\alpha \ll 1$ aggiorno molto lentamente la value function.

If $\alpha = 1/k$ aggiorno la value function in modo da tendere al valore atteso. Devo memorizzare le occorrenze dello stato s.

If $\alpha = \text{cost}$. Aggiorno la value function, pesando maggiormente i risultati collezionati dalle visite dello stato più recenti.



Sommario



Value iteration

Esempi

Learning with temporal differences