



## Sommario



Determinazione ricorsiva della Value function

Determinazione ricorsiva della policy ottima.

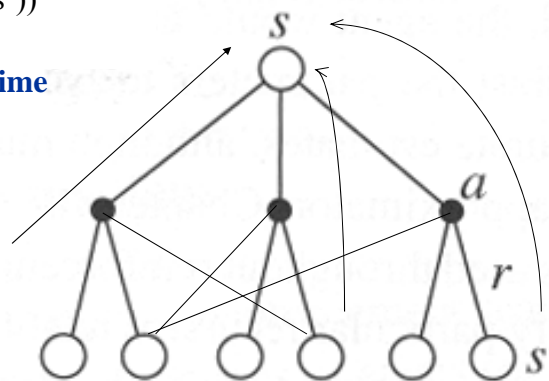


## Osservazioni



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

**Backwards in time  
(full backup)**



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



## Algoritmo per "iterative policy evaluation", versione batch



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

Definiamo una soglia di convergenza  $\tau$

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

3/31

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



## 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).

A.A. 2011-2012

4/31

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



## Esempio



• L'agente evolve esplorando gli stati 1-14. TS sono gli stati terminali.

• L'agente può spostarsi: {dx, sx, alto, basso} se supera il bordo rimane nello stato di partenza.

•  $\pi(s,a)$  è equiprobabile.

•  $R = -1$  per ogni transizione, tranne quando  $s' = TS$  ( $R = 0$ ).

•  $\gamma = 1$ .

TS	1	2	3
4	5	6	7
8	9	10	11
12	13	14	TS



## Esempio - $V(s)_{k=0}$



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



## Esempio - $V(s)_{k=1}$



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$V(s)_{k=0}$

0	-0.75	-1	-1
-0.75	-1	-1	-1
-1	-1	-1	-0.75
-1	-1	-0.75	0

$V(s)_{k=1}$

Esempi:

- $V(s=10)_{k=1} = 0.25 [1 (-1 + 1 V(14))_{k=0}] + 0.25 [1 (-1 + 1 V(6))_{k=0}] + 0.25 [1 (-1 + 1 V(9))_{k=0}] + 0.25 [1 (-1 + 1 V(11))_{k=0}] = -1$
- $V(s=1)_{k=1} = 0.25 [1 (-1 + 1 V(1))_{k=0}] + 0.25 [1 (-1 + 1 V(5))_{k=0}] + 0.25 [1 (0 + 1 V(TS))_{k=0}] + 0.25 [1 (-1 + 1 V(2))_{k=0}] = -0.75$



## Esempio - $V(4)_{k=1}$ per trial



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$V(s)_{k=0}$

0	-0.75	-1	-1
-0,75	??	0	0
0	0	0	0
0	0	-0	0

$V(s)_{k=1}$

Per epoch:

- $V(s=5)_{k=1} = 0.25 [1 (-1 + 1 V(9))_{k=0}] + 0.25 [1 (-1 + 1 V(6))_{k=0}] + 0.25 [1 (-1 + 1 V(4))_{k=0}] + 0.25 [1 (-1 + 1 V(1))_{k=0}] = -1$

Per trial:

- $V(s=5)_{k=1} = 0.25 [1 (-1 + 1 V(9))] + 0.25 [1 (-1 + 1 V(6))] + 0.25 [1 (0 + 1 V(4))] + 0.25 [1 (-1 + 1 V(1))] = -1.375$



## Esempio - $V(s)_{k=2}$



0	-0.75	-1	-1
-0.75	-1	-1	-1
-1	-1	-1	-0.75
-1	-1	-0.75	0

$V(s)_{k=1}$

0	-1.4375	-1.9375	-2
-1.4375	-1.875	-2	-1.9375
-1.9375	-2	-1.875	-1.4375
-2	-1.9375	-1.4375	0

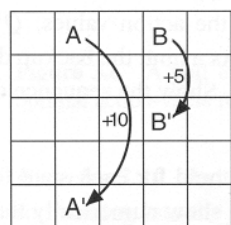
$V(s)_{k=2}$

Esempi:

- $V(s=10)_{k=2} = 0.25 [1 (-1 + 1 V(14)_{k=1})] + 0.25 [1 (-1 + 1 V(6)_{k=1})] + 0.25 [1 (-1 + 1 V(9)_{k=1})] + 0.25 [1 (-1 + 1 V(11)_{k=0})] = -1.875$
- $V(s=1)_{k=2} = 0.25 [1 (-1 + 1 V(1)_{k=1})] + 0.25 [1 (-1 + 1 V(5)_{k=1})] + 0.25 [1 (0 + 1 V(TS)_{k=1})] + 0.25 [1 (-1 + 1 V(2)_{k=1})] = -1.4375$



## Esempio di funzione valore (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)



## Calcolo della funzione valore per uno stato



$t: s = s_t = C$

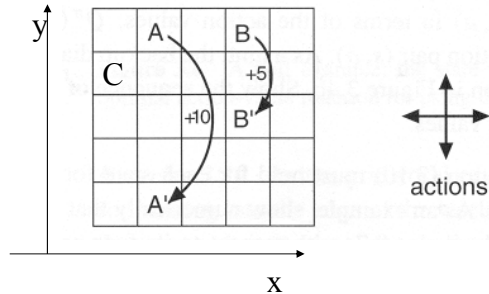
$t + 1: s = s' = s_{t+1}$

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. 2011-2012

11/31

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



## Esempio di valore di uno stato



$t: s = s_t = C$

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

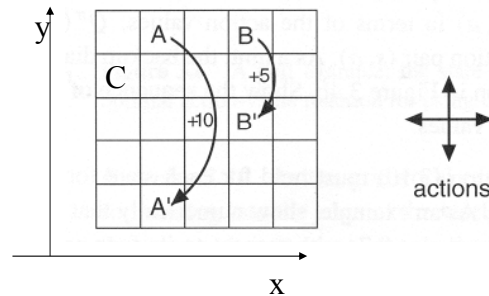
$C + Dy$   $p = 0.25$   $r = 0$

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

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

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

$r_{t+1} = -0.25$



$t + 2, s = C + Dx$

→ A  $p = 0.25$   $r = 0$

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

→  $A + 2Dx$   $p = 0.25$   $r = 0$

→ C  $p = 0.25$   $r = 0$

$r_{t+2} = 0.0$

$s = C - Dy$

→ C  $r = 0$

→  $C - 2Dy$   $r = 0$

→  $C - Dy + Dx$   $r = 0$

→  $C - Dy$   $r = -1$

$r_{t+2} = -0.25$

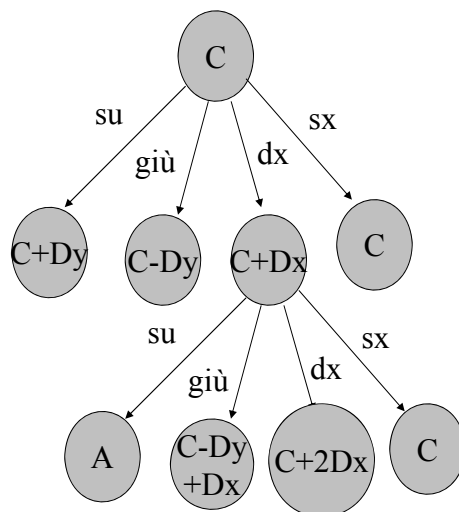
A.A. 2011-2012

12/31

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



## Sviluppo parziale a 2 passi



## Sommario

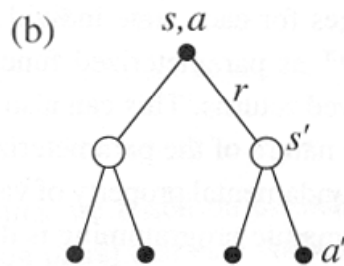
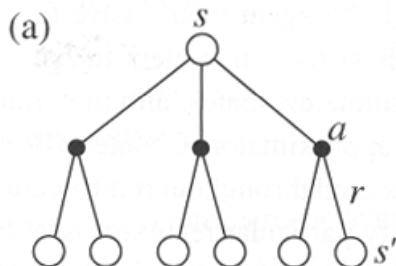


Determinazione ricorsiva della Value function

**Determinazione della policy ottima.**



## 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.

A.A. 2011-2012

15/31

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



## 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,  $\pi_1$ , migliore di un'altra,  $\pi_2$ , se e solo se:

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

In particolare si definisce una politica ottima,  $\pi^*$ , se e solo se:

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

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

A.A. 2011-2012

16/31

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





## 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.



## Esempio: robot



$$V_h = \Pr(W) \times 1 \times [1 + 0.8V_h] + \Pr(S) \times 0.4 \times [3 + 0.8V_h] + \Pr(S) \times 0.6 \times [3 + 0.8V_l]$$

$$V_l = \Pr(W) \times 1 \times [1 + 0.8V_l] + \Pr(S) \times 0.1 \times [3 + 0.8V_l] + \Pr(S) \times 0.9 \times [-3 + 0.8V_h] + \Pr(R) \times 1 \times [0 + 0.8V_h]$$

$$\alpha = 0.8, \beta = 0.4, R_w = 1, R_s = 3, R_r = 0$$

**Sono possibili un numero limitato di scelte deterministiche di azioni**

### STATO: $V_h$

a = search:  $V_h \cong 6 + 0.65V_l$

a = wait:  $V_h = 5$  (per ogni  $V_l$ )

### STATO: $V_l$

a = search:  $V_l \cong -2.2 + 6.6V_h$

a = wait:  $V_l = 5$

a = recharge:  $V_l = 0.8V_h$





## Esempio: robot - policy I



$$V_h = \Pr(W) \times 1 \times [1 + 0.8V_h] + \Pr(S) \times 0.4 \times [3 + 0.8V_h] + \Pr(S) \times 0.6 \times [3 + 0.8V_l]$$

$$V_l = \Pr(W) \times 1 \times [1 + 0.8V_l] + \Pr(S) \times 0.1 \times [3 + 0.8V_l] + \Pr(S) \times 0.9 \times [-3 + 0.8V_h] + \Pr(R) \times 1 \times [0 + 0.8V_h]$$

### Policy iniziale equiprobabile

Value function inizialmente:  $V_h = V_l = 0$

$$V_h \cong 0.5 \times (6 + 0.65V_l) + 0.5 \times 5 = 5.5$$

$$V_l \cong 0.33 \times (-2.2 + 6.6V_h) + 0.33 \times 5 + 0.33 \times 0.8 V_h = -0.726 + 1.65 = 0.924$$



Devo migliorare la stima di  $V(s)$

A.A. 2011-2012

19/31

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



## Esempio: robot - policy I



$$V_h = \Pr(W) \times 1 \times [1 + 0.8V_h] + \Pr(S) \times 0.4 \times [3 + 0.8V_h] + \Pr(S) \times 0.6 \times [3 + 0.8V_l]$$

$$V_l = \Pr(W) \times 1 \times [1 + 0.8V_l] + \Pr(S) \times 0.1 \times [3 + 0.8V_l] + \Pr(S) \times 0.9 \times [-3 + 0.8V_h] + \Pr(R) \times 1 \times [0 + 0.8V_h]$$

### Policy iniziale equiprobabile

Value function equal to:  $V_h = 5.5$

$$V_l = 0.924$$

$$V_h \cong 0.5 \times (6 + 0.65V_l) + 0.5 \times 5 = 5.8$$

$$V_l \cong 0.33 \times (-2.2 + 6.6V_h) + 0.33 \times 5 + 0.33 \times 0.8 V_h = 12.870$$



A.A. 2011-2012

20/31

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



## Esempio: robot - policy I correct value



$$V_h = \Pr(W) \times 1 \times [1 + 0.8V_h] + \Pr(S) \times 0.4 \times [3 + 0.8V_h] + \Pr(S) \times 0.6 \times [3 + 0.8V_l]$$

$$V_l = \Pr(W) \times 1 \times [1 + 0.8V_l] + \Pr(S) \times 0.1 \times [3 + 0.8V_l] + \Pr(S) \times 0.9 \times [-3 + 0.8V_h] + \Pr(R) \times 1 \times [0 + 0.8V_h]$$

**Policy iniziale equiprobabile**

**Value function equal to:**  $V_h = 5.5$

$$V_l = 0.924$$

$$V_h \cong 0.5 \times (6 + 0.65V_l) + 0.5 \times 5 = 5.8$$

$$V_l \cong 0.33 \times (-2.2 + 6.6V_h) + 0.33 \times 5 + 0.33 \times 0.8 V_h = 12.870$$



## Esempio: robot - II



$$V_h = \Pr(W) \times 1 \times [1 + 0.8V_h] + \Pr(S) \times 0.4 \times [3 + 0.8V_h] + \Pr(S) \times 0.6 \times [3 + 0.8V_l]$$

$$V_l = \Pr(W) \times 1 \times [1 + 0.8V_l] + \Pr(S) \times 0.1 \times [3 + 0.8V_l] + \Pr(S) \times 0.9 \times [-3 + 0.8V_h] + \Pr(R) \times 1 \times [0 + 0.8V_h]$$

**Policy iniziale deterministica:**

**STATO:  $V_h$**

$$a = \text{search} \rightarrow V_h \cong 6 + 0.65V_l \cong 9.25$$

**STATO:  $V_l$**

$$a = \text{wait} \rightarrow V_l = 5$$

Posso migliorare la policy?





## Esempio: robot - III



Miglioro la policy, modificando l'azione associata a V1:

**STATO: Vh**

a = search  $\rightarrow V_h \cong 6 + 0.65V_1 \cong 9.25$

**STATO: V1**

a = recharge  $\rightarrow V_1 = 0.8V_h = 7,4$

Ho stimato correttamente V(s)? No



**STATO: Vh**

a = search  $\rightarrow V_h \cong 6 + 0.65V_1 \cong 11.06$

**STATO: V1**

a = recharge  $\rightarrow V_1 = 0.8V_h = 7,4$

Ho stimato correttamente V(s)? No. Devo iterare la policy evaluation.



## Esempio: robot - IV



Asintoticamente calcolo il valore degli stati:

**STATO: Vh**

a = search  $\rightarrow V_h \cong 6 + 0.65V_1 \rightarrow 12.5$

**STATO: V1**

a = recharge  $\rightarrow V_1 = 0.8V_h \rightarrow 10.0$



Potrei ottenere gli stessi valori risolvendo il sistema lineare:

$$V_h = 6 + 0.65V_1 = 12.5$$

$$V_1 = 0.8V_h = 10.0$$

Ho terminato?



## 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-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-stable*  $\leftarrow$  false

If *policy-stable*, then stop; else go to 2

A.A. 2011-2012

25/31

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



## Esempio



Consideriamo uno scommettitore che scommette su testa e croce.

Se esce testa, vince tanti Euro quanti ne ha scommesso. Se esce croce, perde gli Euro che ha scommesso.

Il gioco termina quando lo scommettitore arriva ad accumulare 100 Euro di vincita o perde tutto il suo capitale.

Il suo capitale di partenza è variabile tra 1 e 99 Euro. Ad ogni lancio della moneta, lo scommettitore può decidere quanto scommettere sul fatto che esca testa.

Undiscounted, episodico, MDP.

Lo stato è il capitale accumulato dal giocatore:  $s = \{1, 2, \dots, 99\}$ .

L'azione è quanto viene scommesso ad ogni lancio:  $a = \{1, 2, \dots, \min(s, 100 - s)\}$ .

Conosciamo la probabilità  $p$  con cui esce testa (la moneta può essere truccata e quindi le due condizioni testa/croce possono essere non equiprobabili).

La dinamica dell'ambiente non dipende da  $a$  ed è rappresentata dalla transizione di uno stato all'altro che a sua volta è funzione del fatto che esca testa o meno.

Il reward istantaneo = 0 tranne quando raggiunge 100 Euro, nel qual caso reward = +1.

Vogliamo massimizzare la probabilità di vincere. Si può fare tramite Value iteration.

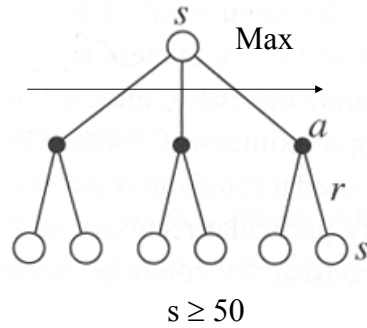
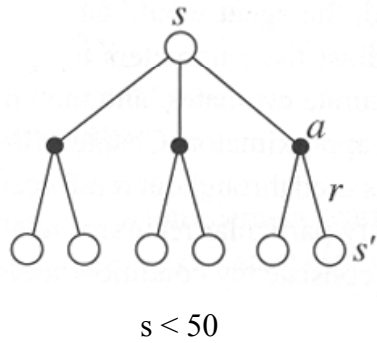
A..

se\



# 1° passo di Value iteration

$P_{s \rightarrow s'} = \{\text{testa}, \text{croce}\}$  non dipende da  $a$ . Sono indipendenti.



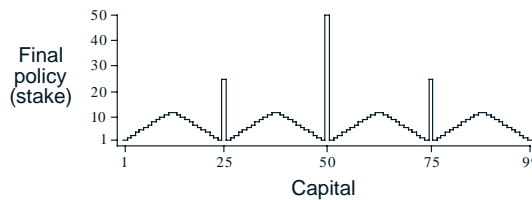
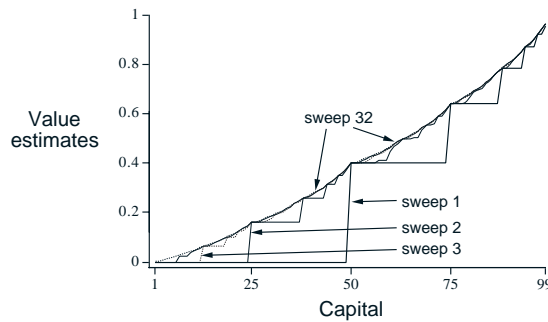
$$V_{k+1}(s) = 0.4 (0 + \gamma V(s')) = 0 \quad \forall a$$

$$V_{k+1}(s) = 0.4 (1 + \gamma V(s')) = 0.4$$

Per  $a$  scelto in modo da arrivare a 100€ in caso di vincita.

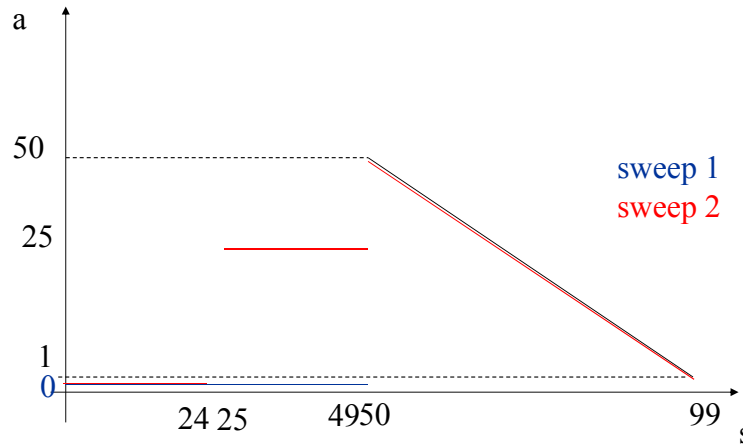


# Stima della Value function





## Policy selezionata



Ci sono diversi stati per cui la  $Q(s, \cdot)$  assume lo stesso valore per diverse azioni.

A.A. 2011-2012

29/31

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



## Esercizio: autonoleggio



2 locazioni di autonoleggio

In ogni locazione, se quando arriva un cliente l'auto è disponibile, si guadagnano 10 euro.

Se l'auto non è disponibile si perde la gestione della locazione.

Le auto diventano disponibili il giorno dopo che sono state restituite dopo il noleggio.

Si possono portare le auto da un autonoleggio all'altro di notte al costo di 2 Euro per ogni auto.

Supponiamo che il numero di auto richieste e restituite in ognuno dei 2 autonoleggi sia rappresentato da una distribuzione di Poisson (probabilità che vengano richieste  $n$  auto:  $(\lambda^n / n!)e^{-\lambda}$  dove  $\lambda$  è il valore atteso (media) di auto richieste o restituite).

Supponiamo che non ci possano essere più di 20 auto in uno dei 2 autonoleggi (ciascuna auto aggiuntiva viene inviata al centro di raccolta nazionale della compagnia).

Supponiamo anche che un massimo di 5 auto possa essere spostato in una singola notte.

Questo problema si può formulare con un MDP dove lo stato è rappresentato dal numero di auto presenti in ciascuno dei 2 autonoleggi al termine di una giornata e le azioni il numero di auto che vengono spostate durante la notte.

Consideriamo  $\gamma = 0.9$ . Il reward sarà il reward accumulato durante il giorno + notte.

**Partiamo dalla policy  $\pi(s, a) = 0$ : nessuna auto viene mossa. Siamo in grado di migliorarla? Come?**

A.A. 2011-2012

30/31

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



## Sommario



Determinazione ricorsiva della Value function

Determinazione della policy ottima.