

# **PROCESSOS DE DECISÃO SEQUENCIAL**

Luís Morgado

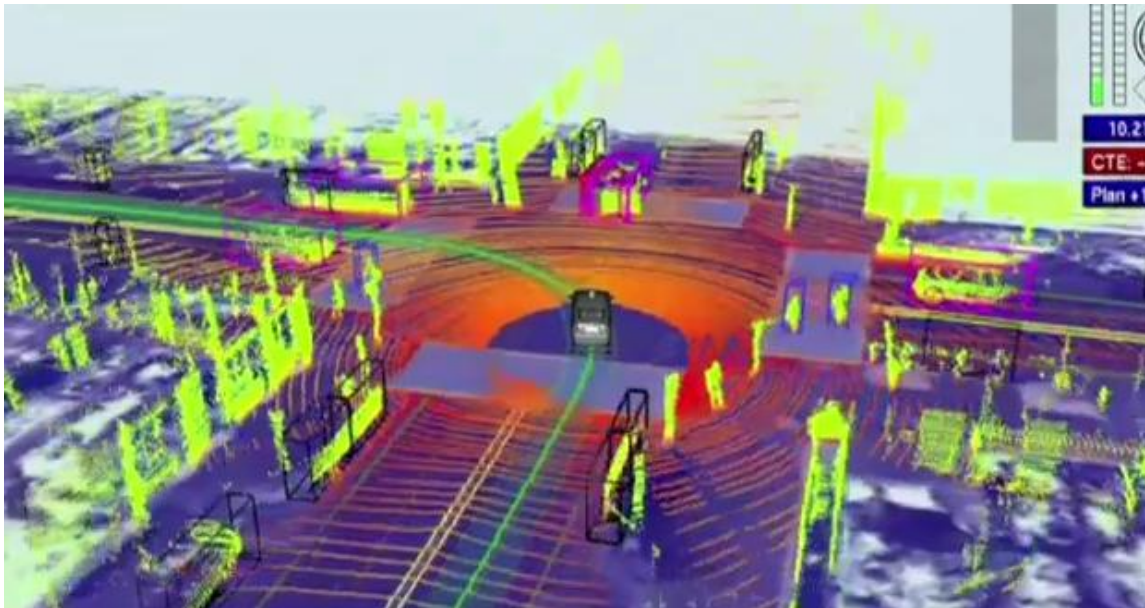
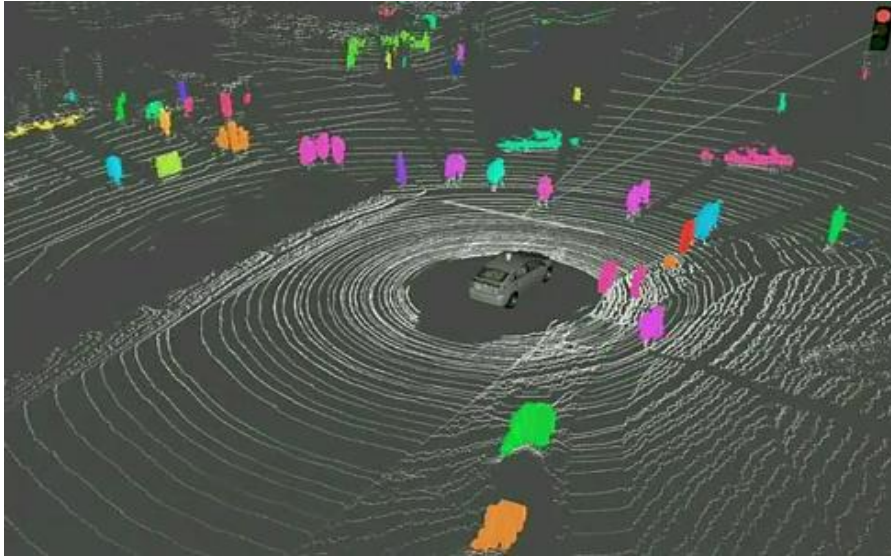
2015

# Processos de Decisão Sequencial

---

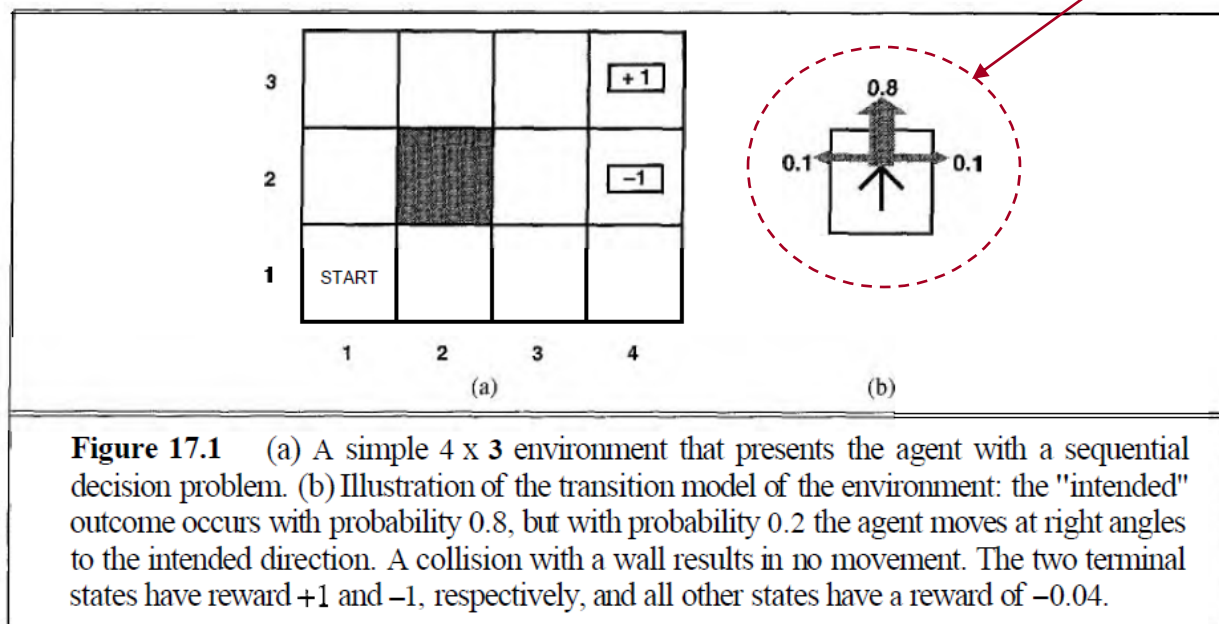
- Muito estudados desde os anos 50
- Controlo de Sistemas
- Investigação Operacional
  - Planeamento, escalonamento, logística
  - Economia, mercados financeiros, telecomunicações
- Inteligência artificial (final  $\approx 80$ )
  - Aprendizagem por reforço
  - Planeamento probabilístico
- Programação Dinâmica o principal método utilizado

# Processos de Decisão Sequencial

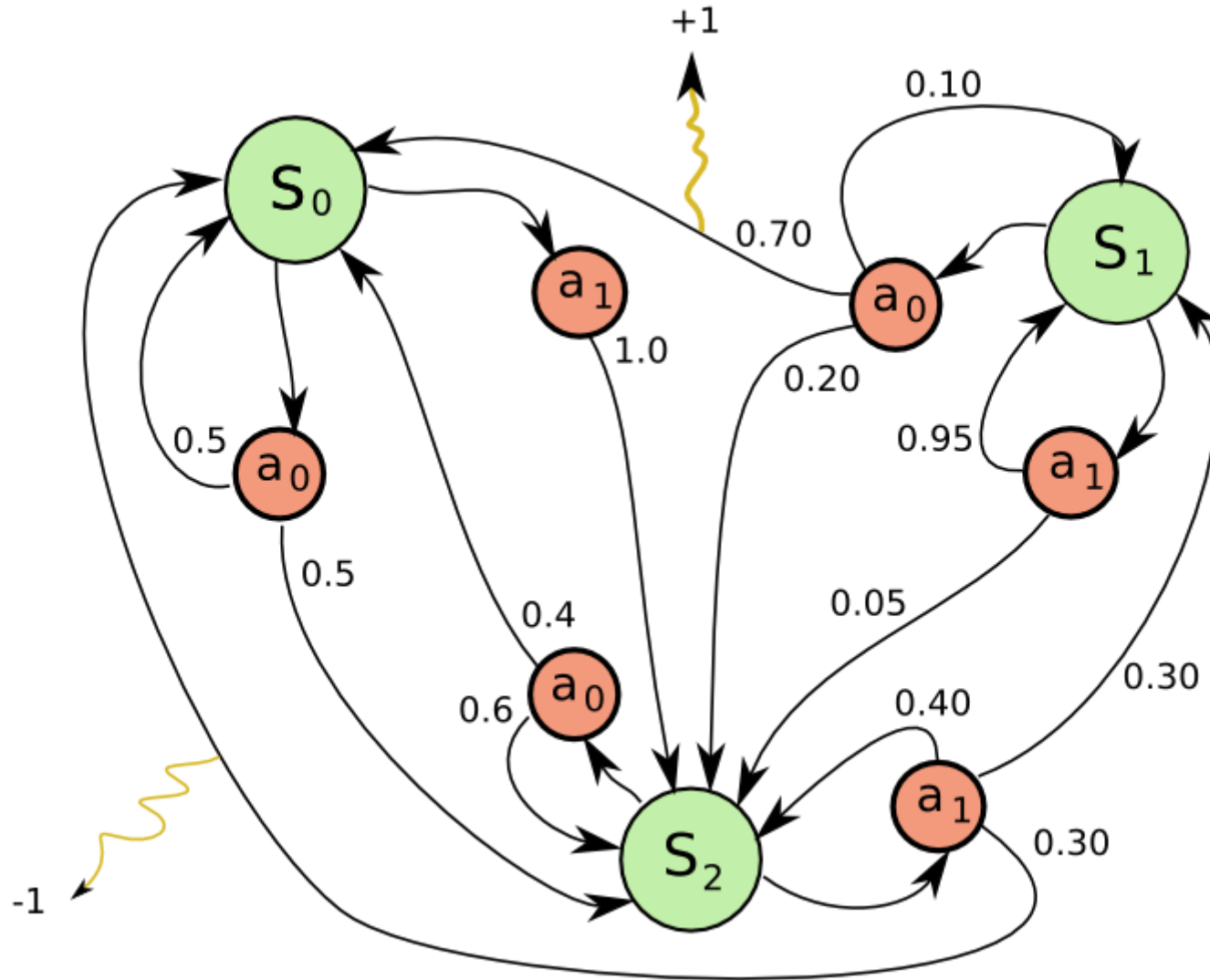


# Processos de Decisão Sequencial

- Problema da decisão ao longo do tempo
  - Utilidade de uma acção depende de uma sequência de decisões
  - Possibilidade de ganhos e perdas
  - Incerteza na decisão
  - Efeito cumulativo



# Processos de Decisão Sequencial



# Propriedade de Markov

---

- Andrey Markov
  - Matemático Russo (1856 – 1922)
- Um processo estocástico tem a ***propriedade de Markov*** se a distribuição probabilística condicional dos **estados futuros** de um processo depender exclusivamente do **estado presente**
- **A previsão dos estados seguintes só depende do estado presente**

# Processos de Decisão de Markov

---

- Representação do mundo sob a forma de PDM

$S$  – conjunto de estados do mundo

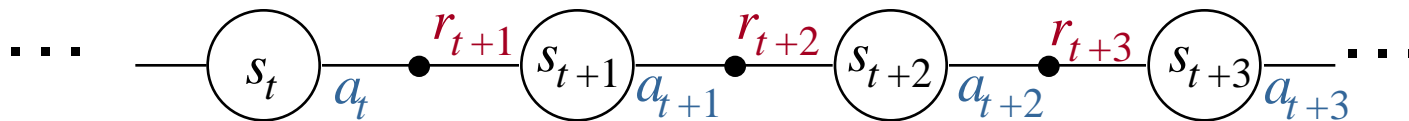
$A(s)$  – conjunto de acções possíveis no estado  $s \in S$

$T(s, a, s')$  – probabilidade de transição de  $s$  para  $s'$  através de  $a$

$R(s, a, s')$  – retorno esperado na transição de  $s$  para  $s'$  através de  $a$

$\gamma$  – taxa de desconto para recompensas diferidas no tempo

$t = 0, 1, 2, \dots$  – tempo discreto



**Cadeia de Markov**

# Utilidade

---

Efeito cumulativo da evolução da situação

- História de evolução  $h$ 
  - Sequência de estados (com ganhos/perdas)
- Recompensa
  - Ganho ou perda num determinado estado
  - Valor finito positivo ou negativo
  - $R(s)$
- $U_h([s_0, s_1, \dots, s_n])$



# Utilidade

---

- Recompensas aditivas

- $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$

- Recompensas descontadas (no tempo)

- $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$

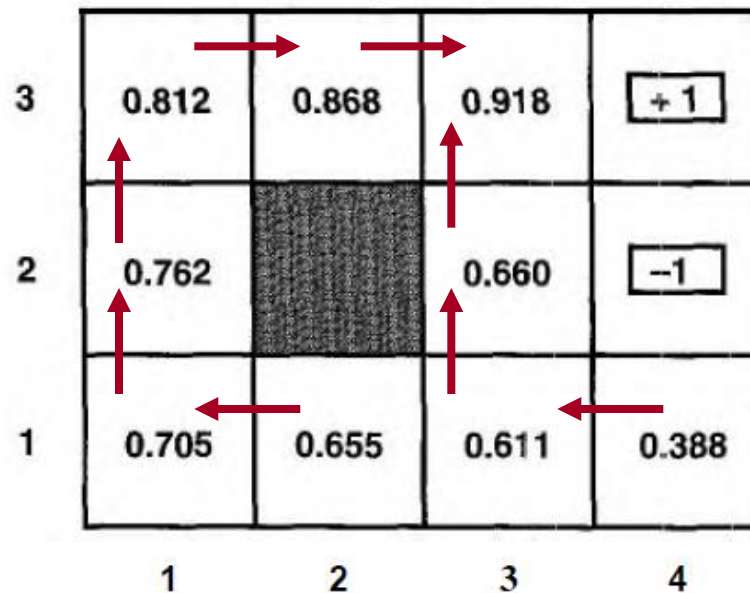
- Factor de desconto

- $\gamma \in [0,1]$

- Recompensas não estão limitadas a uma gama finita de valores

# Utilidade (valor) de estado

Exemplo:



**Figure 17.3** The utilities of the states in the 4 x 3 world, calculated with  $\gamma = 1$  and  $R(s) = -0.04$  for nonterminal states.

# Política Comportamental

---

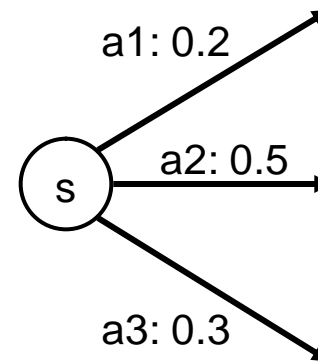
- Forma de representação do comportamento do agente
- Define qual a acção que deve ser realizada em cada estado (estratégia de acção)

- Política **determinista**

$$\pi : S \rightarrow A(s) ; s \in S$$

- Política **não determinista**

$$\pi : S \times A(s) \rightarrow [0,1] ; s \in S$$



# O Princípio da Solução Óptima

---

- Programação Dinâmica
  - Requer a decomposição em sub-problemas
- Num PDM isso deriva da assunção da independência dos caminhos
- As utilidades dos estados podem ser determinados em função das utilidades dos estados sucessores

$$U^{\pi}(s) = E\langle r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \rangle$$

$$= E\langle r_1 + \gamma U^{\pi}(s') \rangle$$

$$= \sum_a \pi(s, a) \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^{\pi}(s')]$$

**Equações de  
Bellman**

# Processos de Decisão de Markov

---

Utilidade de estado para uma política  $\pi$

$$U^{\pi}(s) = \sum_a \pi(s, a) \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^{\pi}(s')]$$

Política óptima  $\pi^*$

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U(s')]$$

Utilidade de estado para a política óptima  $\pi^*$

$$U^{\pi^*}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U^{\pi^*}(s')]$$

# Processos de Decisão de Markov

---

## Iteração da utilidade de estado

Iniciar  $U(s)$ :

$$U(s) \leftarrow 0, \quad \forall s \in S$$

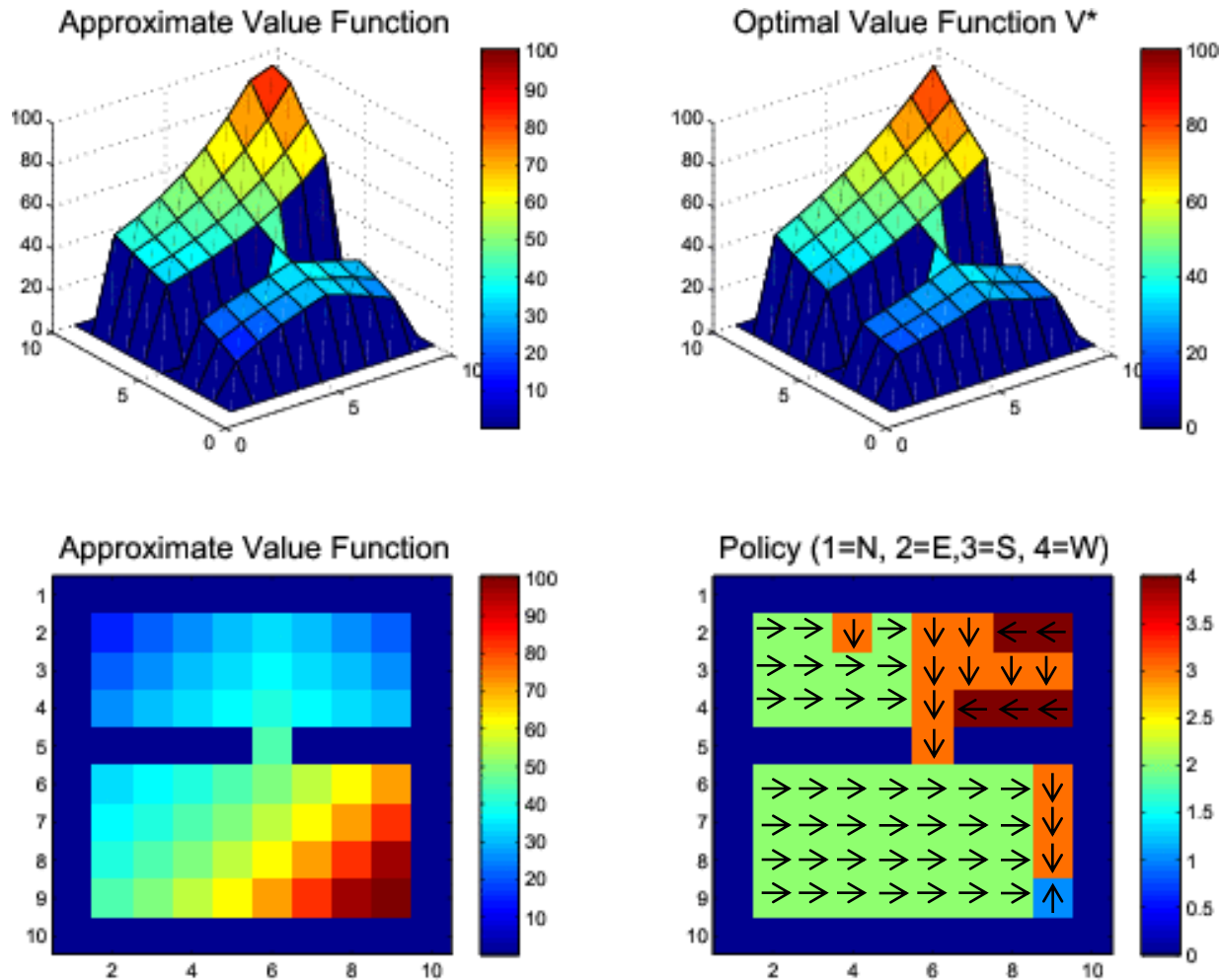
Iterar  $U(s)$ :

$$U(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U(s')], \quad \forall s \in S$$

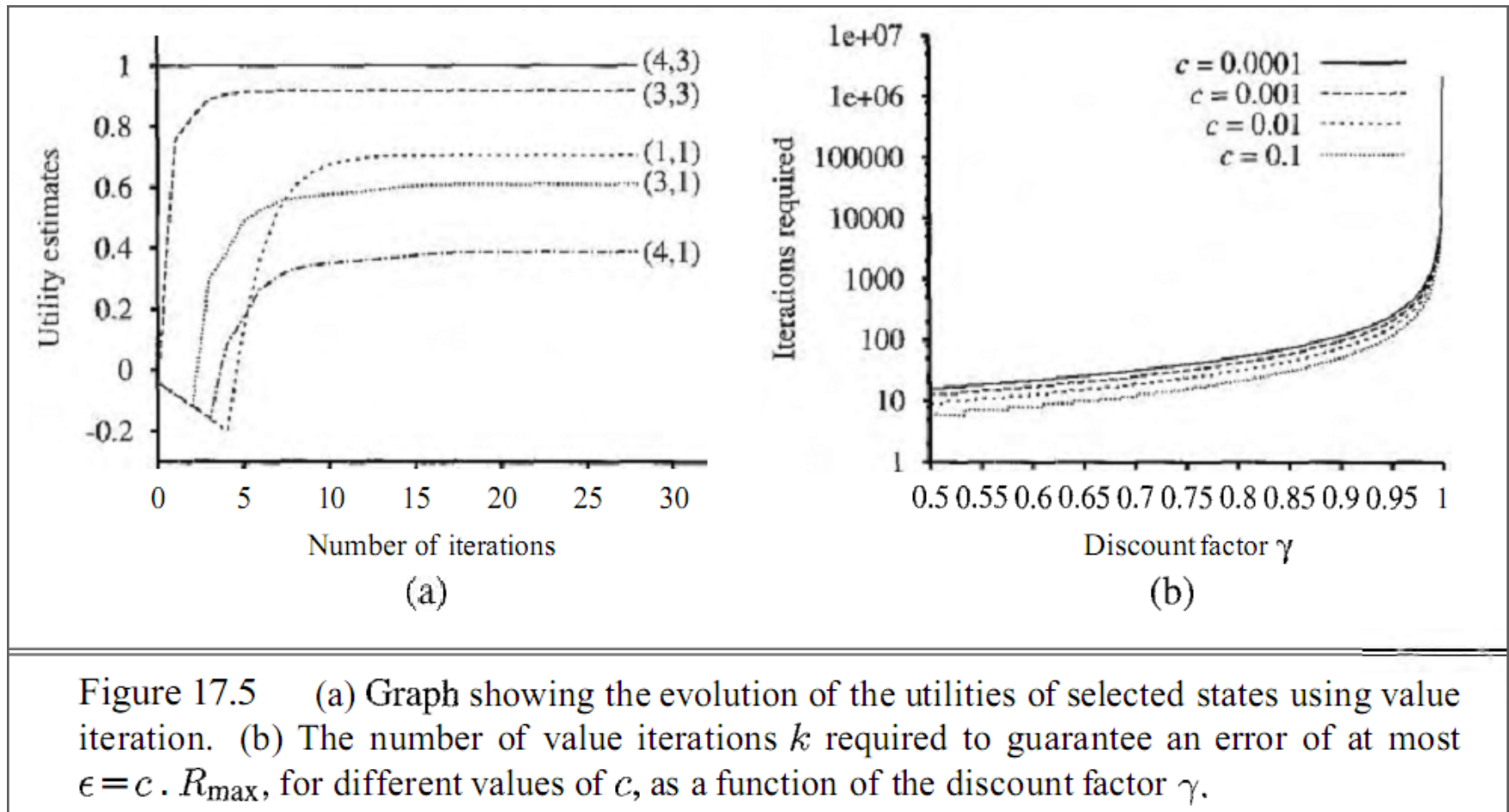
No limite:

$$U \rightarrow U^{\pi^*}$$

# Processos de Decisão de Markov

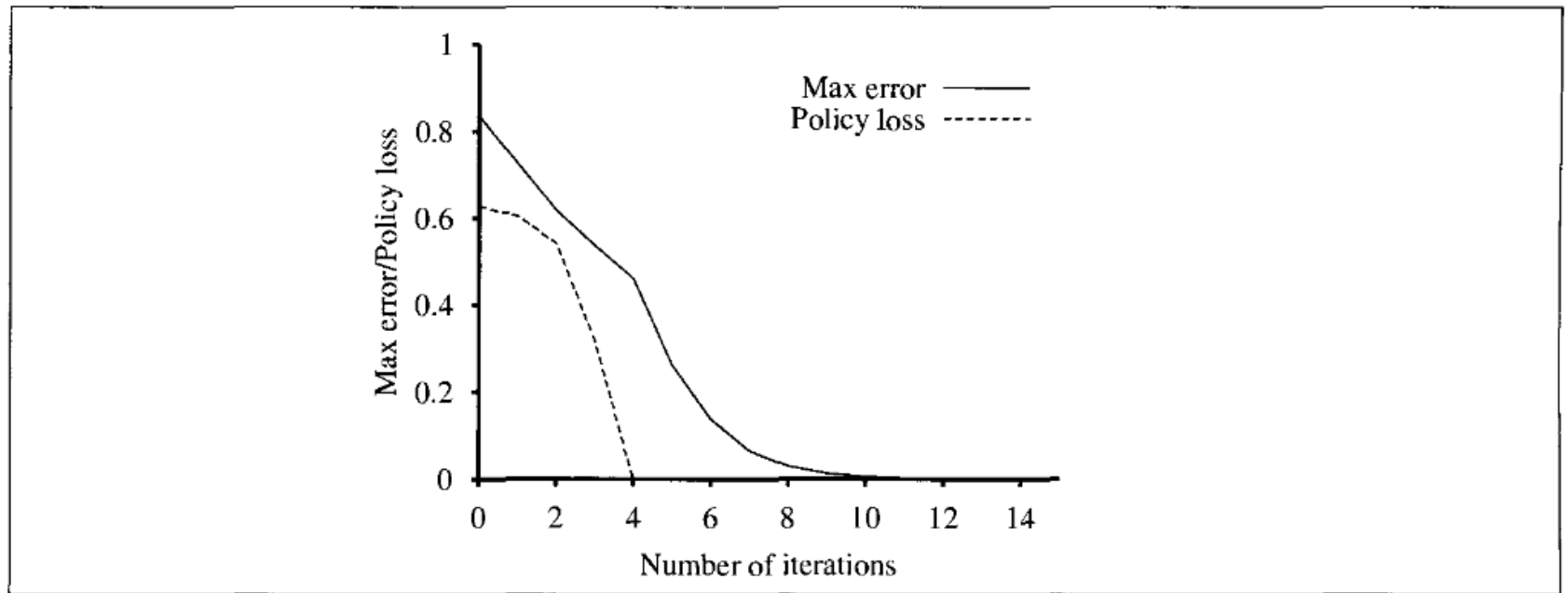


# Cálculo da Utilidade de Estado





# Cálculo da Utilidade de Estado



**Figure 17.6** The maximum error  $\|U_i - U\|$  of the utility estimates and the policy loss  $\|U^{\pi_i} - U\|$  compared with the optimal policy, as a function of the number of iterations of value iteration.

if  $\|U_{i+1} - U_i\| < \epsilon(1 - \gamma)/\gamma$  then  $\|U_{i+1} - U\| < \epsilon$

# Processos de Decisão de Markov

---

## Iteração da utilidade de estado

Iniciar  $U(s)$ :

$$U(s) \leftarrow 0, \quad \forall s \in S$$

Iterar  $U(s)$ :

$$U(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U(s')], \quad \forall s \in S$$

No limite:

$$U \rightarrow U^{\pi^*}$$

## Critério de paragem de iteração?

- Diferença máxima de actualização  $< \Delta_{\max}$  (limiar de convergência)

# Processos de Decisão de Markov

---

- **Propriedade de Markov**
  - Estados futuros dependem apenas do estado actual
    - São independentes de estados passados
- **Modelo do mundo - representação do problema**
  - Conjunto de estados
    - $S$
  - Conjunto de acções possíveis num estado
    - $A(s)$
  - Modelo de transição
    - $T(s,a,s')$  – também designado  $P(s,a,s')$
  - Modelo de recompensa
    - $R(s,a,s')$  – no caso geral
    - $R(s, a)$  – se a recompensa só depende do estado e da acção
    - $R(s)$  – se a recompensa só depende do estado

# Processos de Decisão de Markov

---

No caso geral

$$U(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U(s')]$$

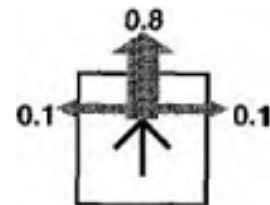
Se a recompensa só depende do estado

$$U(s) = \max_a \sum_{s'} T(s, a, s') [R(s) + \gamma U(s')]$$

$$U(s) = R(s) + \max_a \sum_{s'} T(s, a, s') [\gamma U(s')]$$

# Cálculo da Utilidade de Estado

3	0.812	0.868	0.918	$\boxed{+1}$
2	0.762		0.660	$\boxed{-1}$
1	0.705	0.655	0.611	0.388
	1	2	3	4



$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U_i(s')$$

$$U(1,1) = -0.04 + \gamma \max \left\{ \begin{array}{ll} 0.8U(1,2) + 0.1U(2,1) + 0.1U(1,1), & (Up) \\ 0.9U(1,1) + 0.1U(1,2), & (Left) \\ 0.9U(1,1) + 0.1U(2,1), & (Down) \\ 0.8U(2,1) + 0.1U(1,2) + 0.1U(1,1) \} & (Right) \end{array} \right.$$

We can think of the value iteration algorithm as *propagating information* through the state space by means of local updates.

# Processos de Decisão de Markov

---

- Problemas
  - Dimensão dos espaços de estados  
(“*The curse of dimensionality*”)
  - Dificuldade de definição das dinâmicas  
(por exemplo a partir de dados experimentais)
  - Dinâmicas desconhecidas

$$U(s) \leftarrow \max_{a \in A(s)} \sum_{s'} T(\underset{?}{s}, \underset{?}{a}, s') [R(s, a, s') + \gamma U^\pi(s')] \quad \forall s \in S$$

# Referências

---

[Russel & Norvig, 2003]

S. Russell, P. Norvig, “Artificial Intelligence: A Modern Approach”, 2nd Ed., Prentice Hall, 2003

[Sutton & Barto, 1998]

R. Sutton, A. Barto, “Reinforcement Learning: An Introduction”, MIT Press, 1998

[Mahadevan, 2009]

S. Mahadevan, “Learning Representation and Control in Markov Decision Processes: New Frontiers”, Foundations and Trends in Machine Learning, 1:4, 2009

[LaValle, 2006]

S. LaValle, “Planning Algorithms”, Cambridge University Press, 2006

[Kragic & Vincze, 2009]

D. Kragic, M. Vincze, “Vision for Robotics”, Foundations and Trends in Robotics, 1:1, 2009