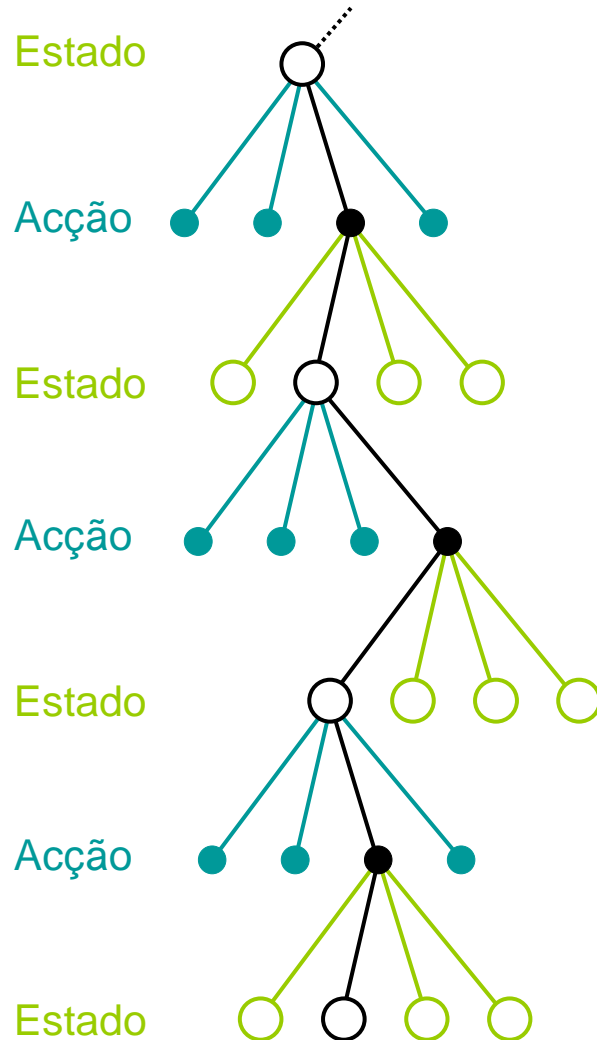


PROCESSOS DE DECISÃO SEQUENCIAL

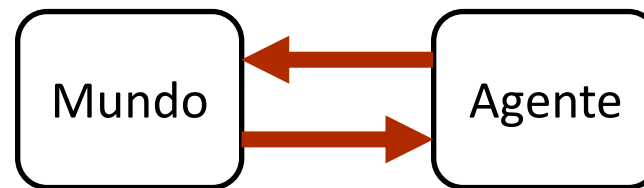
Luís Morgado

2021

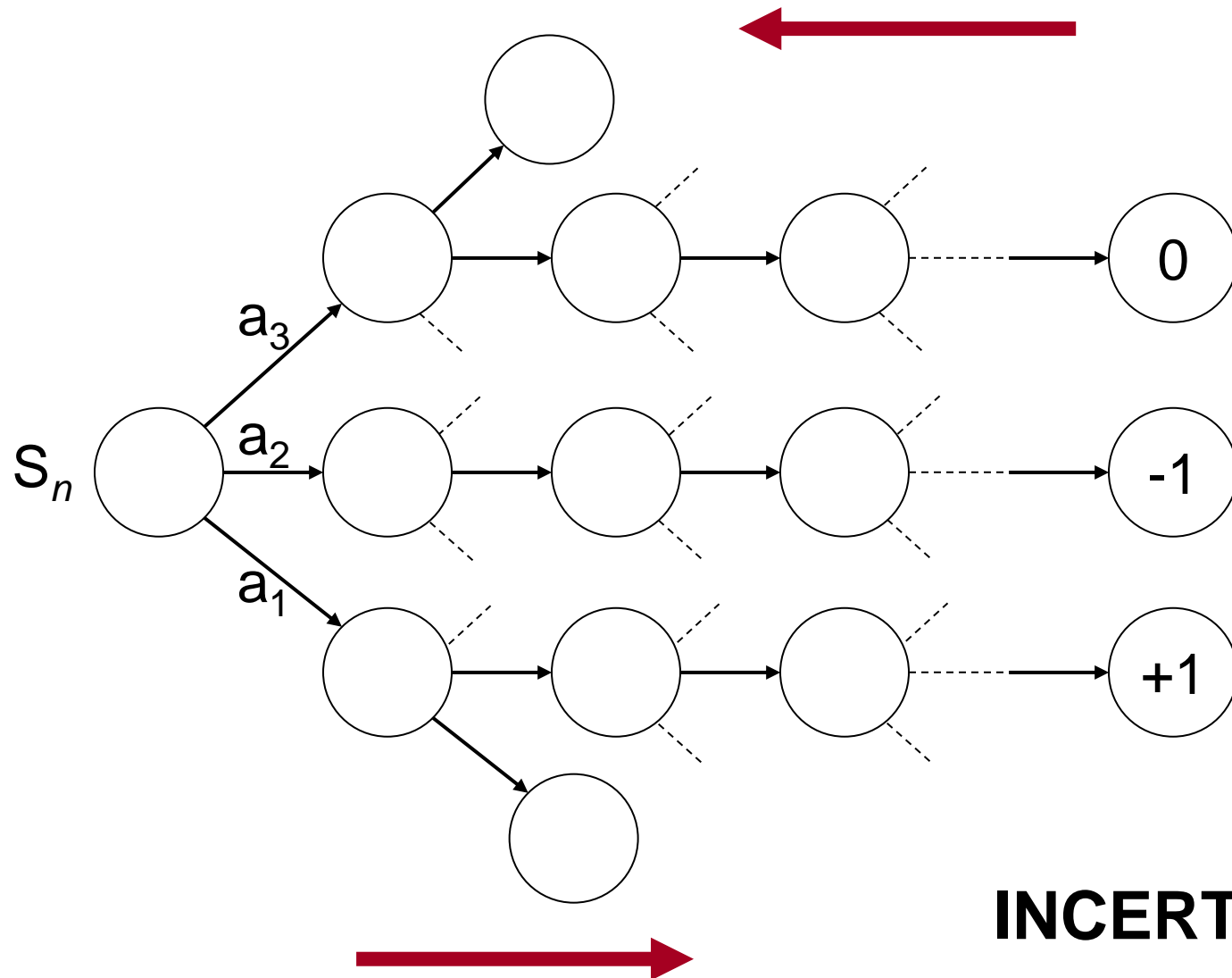
PROCESSOS DE DECISÃO SEQUENCIAL



Como prever e controlar o desenrolar da interacção entre agente e ambiente ao longo do tempo para um **objectivo de longo prazo**?

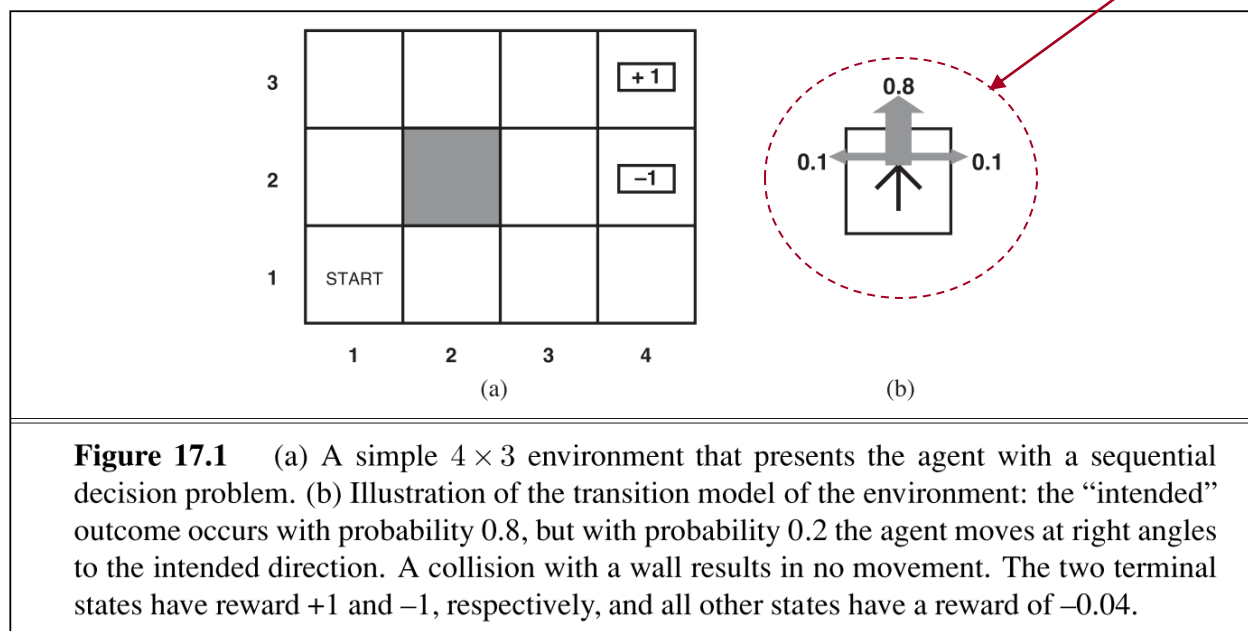


PROCESSOS DE DECISÃO SEQUENCIAL

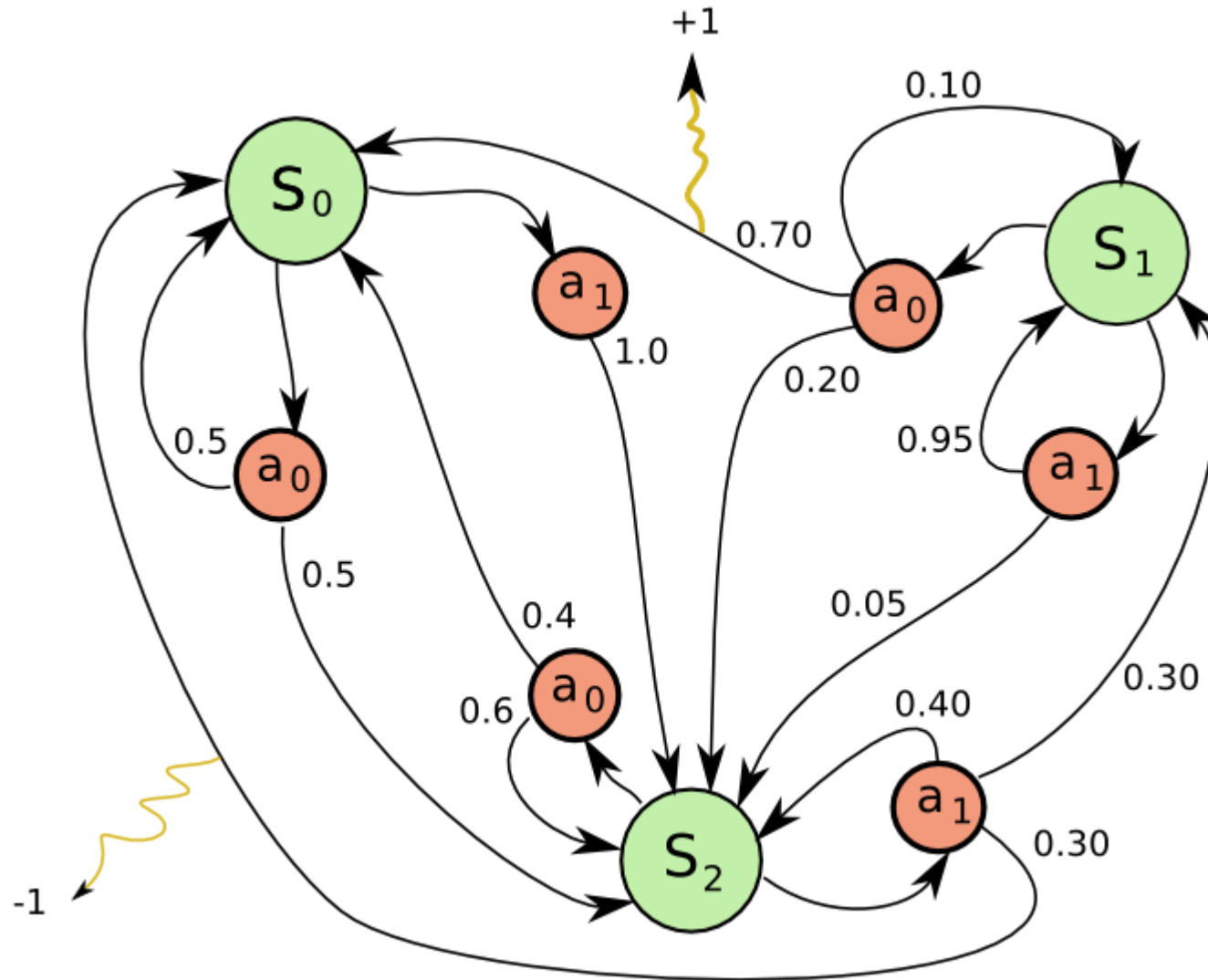


Processos de Decisão Sequencial

- 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

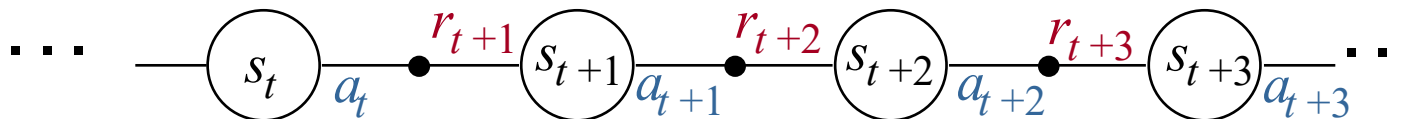
$A(s)$ – conjunto de acções possíveis no estado $s \in S$

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

$R(s, a, s')$ – recompensa esperada na transição de s para s' através de a

γ – 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)

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

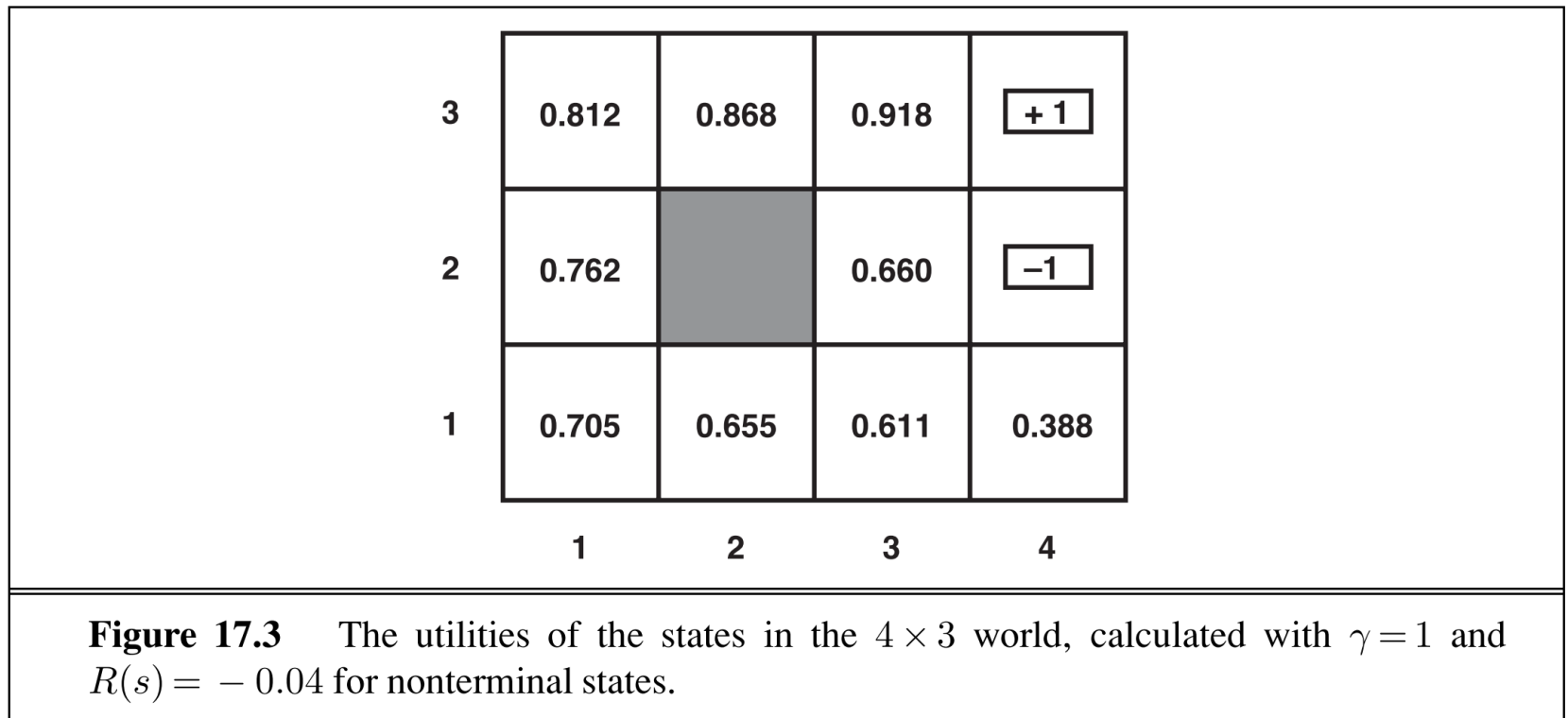
- Factor de desconto

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

- Recompensas não estão limitadas a uma gama finita de valores

Utilidade (valor) de estado

Exemplo:



Política Comportamental

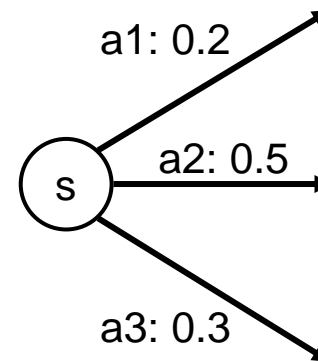
- Forma de representação do comportamento do agente
- Define qual a acção que deve ser realizada em cada estado (estratégia de acção)

- Política **determinista**

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

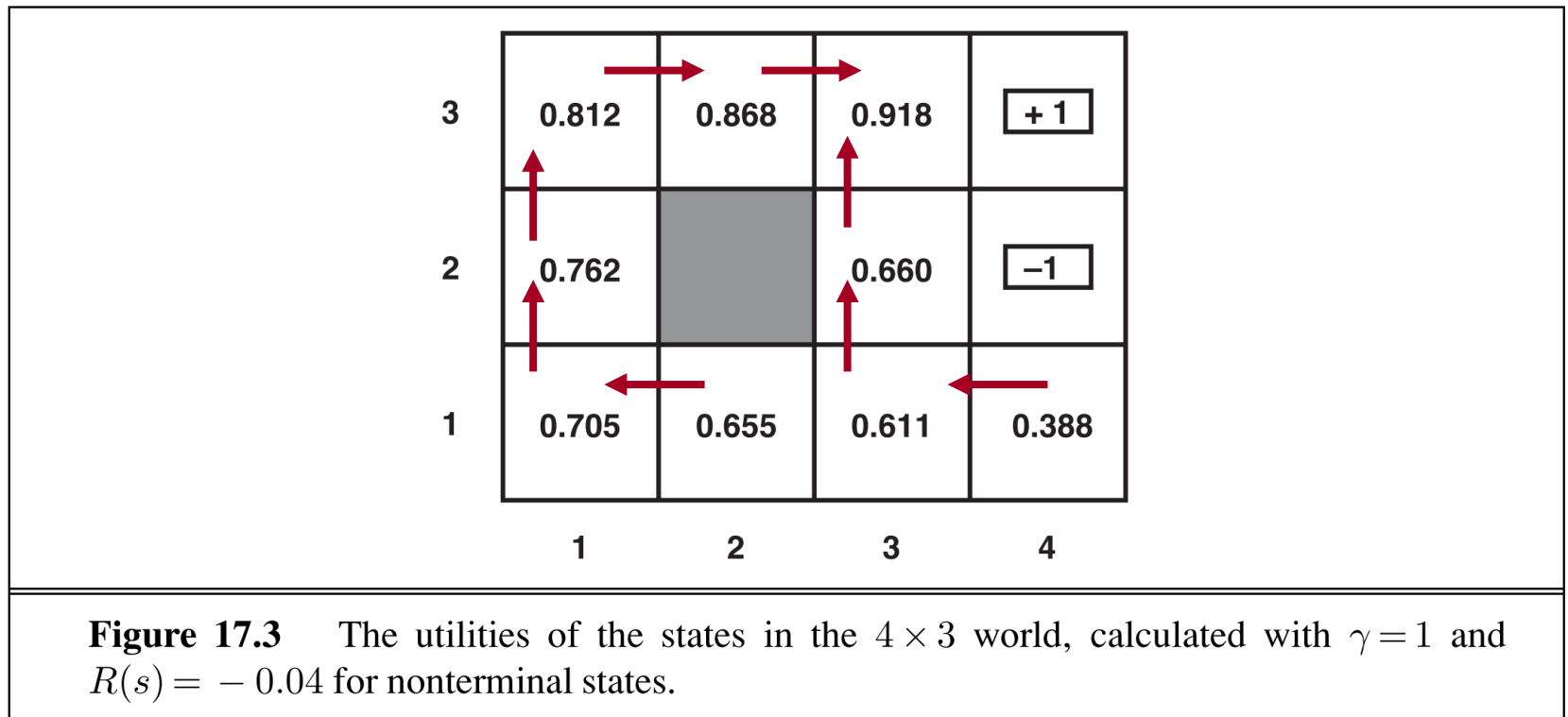
- Política **não determinista**

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



Política Comportamental

Exemplo:



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**

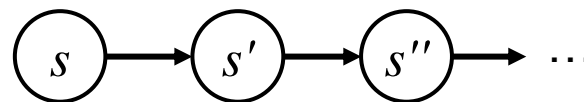
Valor esperado

$$E\langle X \rangle = \sum_{i=0}^{\infty} x_i p(x_i)$$

$$E\langle X + Y \rangle = E\langle X \rangle + E\langle Y \rangle$$

Cadeia de Markov

Política: π

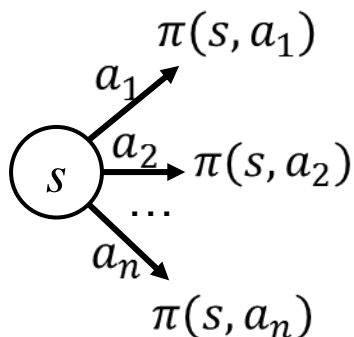


Episódio

1	r_1^1	r_2^1	r_3^1	...
2	r_1^2	r_2^2	r_3^2	...
...				

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

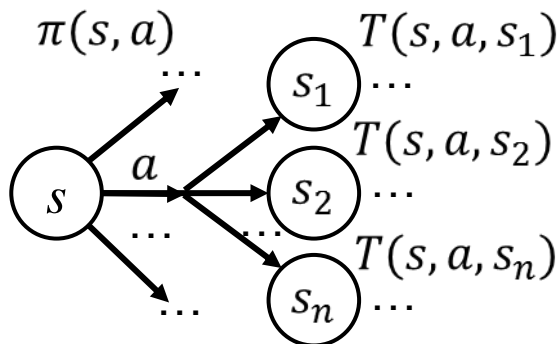
Política



$$\pi(s, a)$$

$$a \in A(s) = \{a_1, a_2, \dots, a_n\}$$

Transição de estado com base num modelo



$$T(s, a, s')$$

$$s' \in \text{suc}(s) \subseteq S$$

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

Processos de Decisão de Markov

Utilidade de estado para uma política π

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

Política óptima π^*

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

Utilidade de estado para a política óptima π^*

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

Processos de Decisão de Markov

Iteração da utilidade de estado

Iniciar $U(s)$:

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

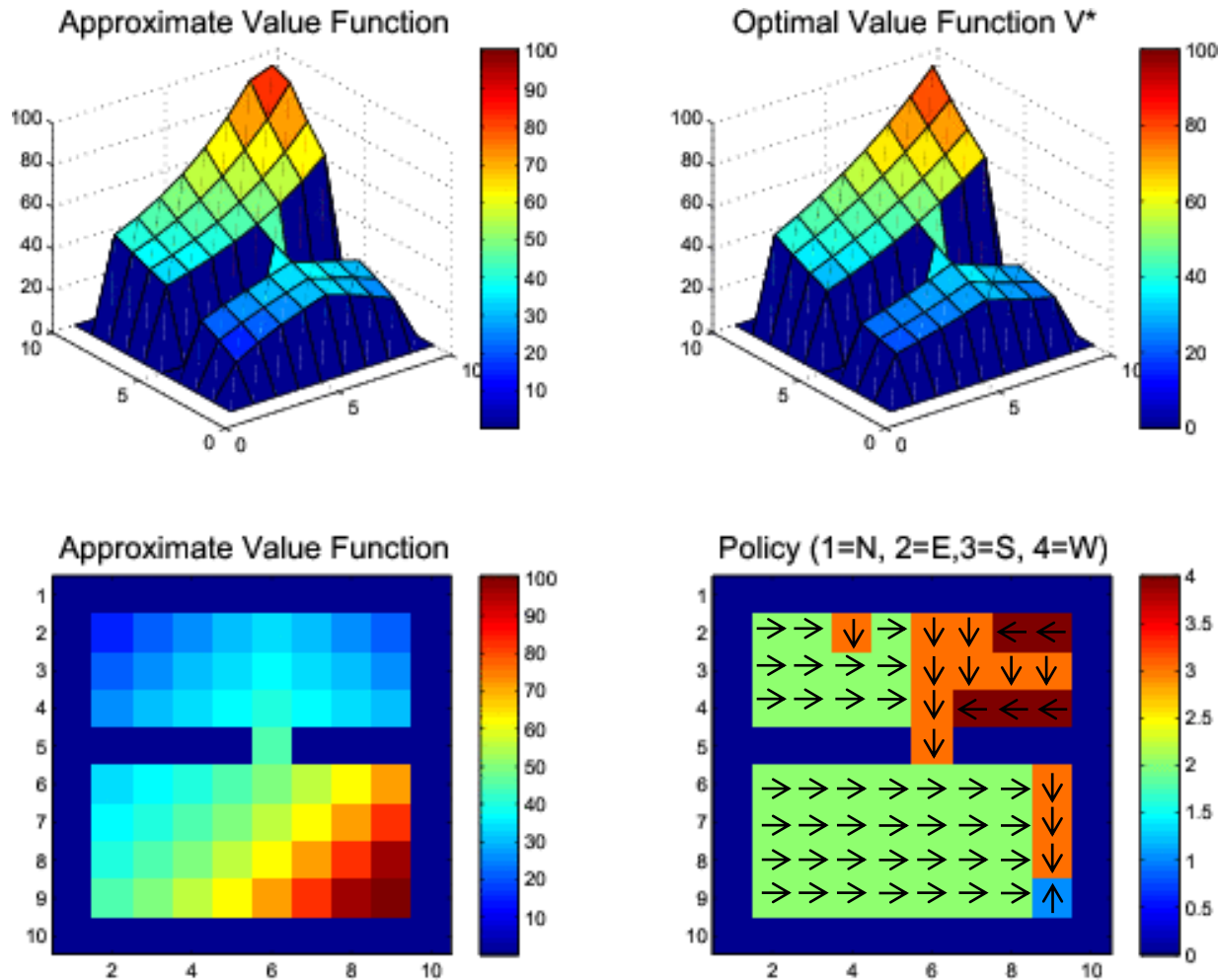
Iterar $U(s)$:

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

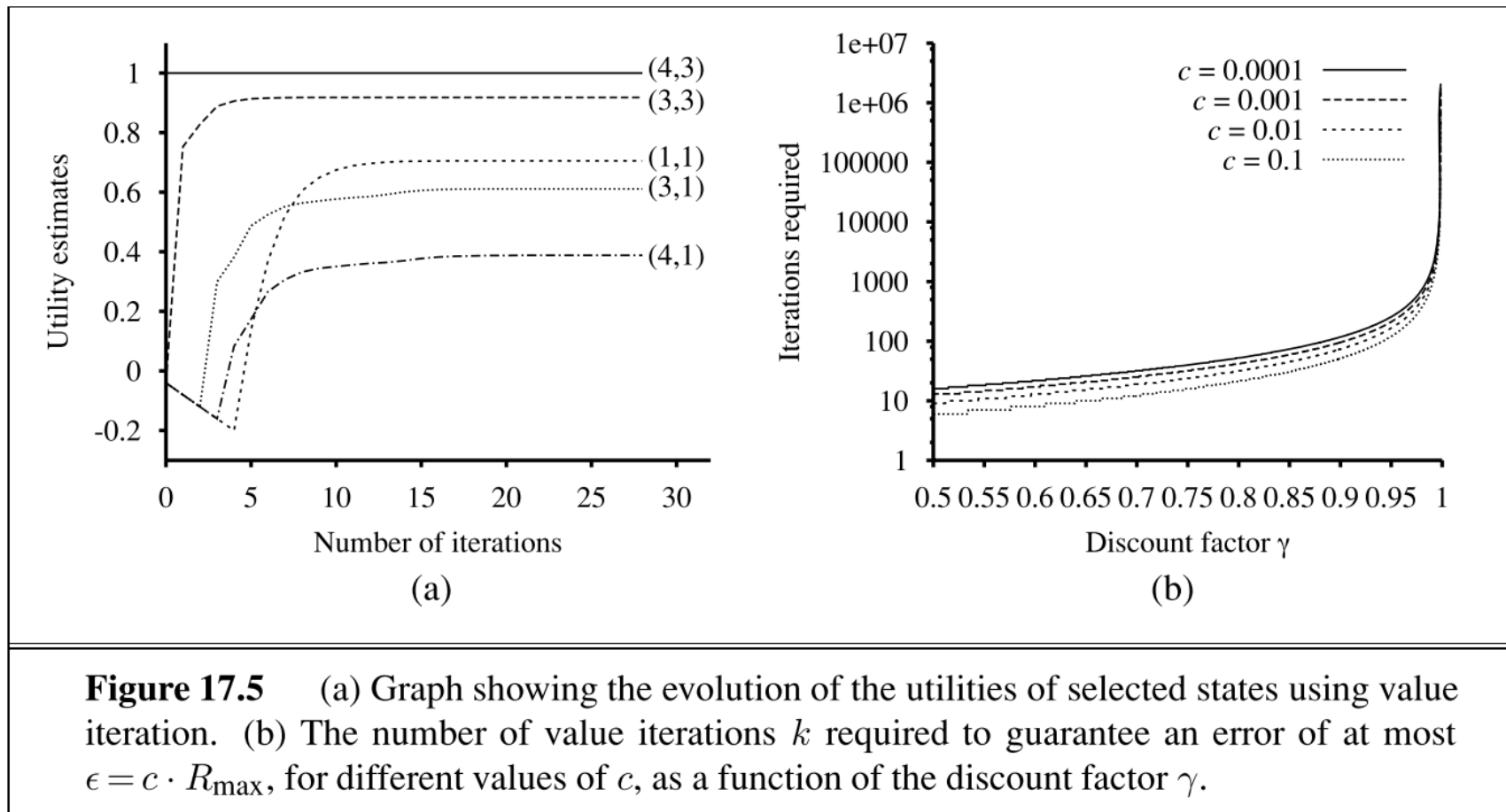
No limite:

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

Processos de Decisão de Markov



Cálculo da Utilidade de Estado



Cálculo da Utilidade de Estado

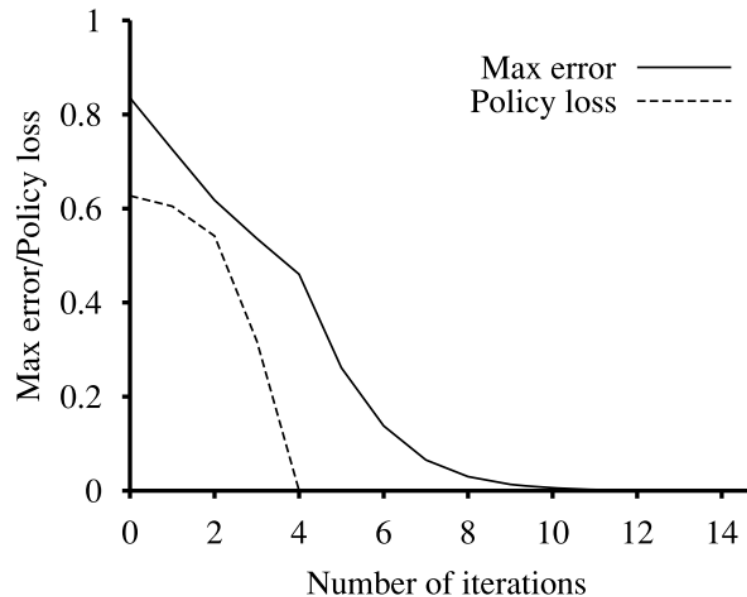


Figure 17.6 The maximum error $\|U_i - U\|$ of the utility estimates and the policy loss $\|U^{\pi_i} - U\|$, as a function of the number of iterations of value iteration.

if $\|U_i - U\| < \epsilon$ then $\|U^{\pi_i} - U\| < 2\epsilon\gamma/(1 - \gamma)$

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 $\leq \Delta_{\max}$ (limiar de convergência)

CÁLCULO DA UTILIDADE

Iteração da utilidade de estado

Iniciar $U(s)$:

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

Iterar $U(s)$:

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

No limite:

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

Critério de paragem de iteração?

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

```
function utilidade:
```

```
     $U \leftarrow 0, \forall s \in S$ 
```

```
  do:
```

```
     $U_{ant} \leftarrow U$ 
```

```
     $\delta \leftarrow 0$ 
```

```
    for  $s$  in  $S$ :
```

```
         $U[s] \leftarrow \max_{a \in A(s)} U_{acção}(s, a, U_{ant})$ 
```

```
         $\delta \leftarrow \max\{\delta, |U[s] - U_{ant}[s]|\}$ 
```

```
  while  $\delta > \Delta_{max}$ :
```

```
  return  $U$ 
```

```
function  $U_{acção}(s, a, U)$ :
```

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

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

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

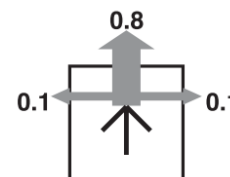
Se a recompensa só depende do estado

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

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

Cálculo da Utilidade de Estado

3	0.812	0.868	0.918	+1
2	0.762		0.660	-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.

Referências

[Russel & Norvig, 2010]

S. Russell, P. Norvig, “Artificial Intelligence: A Modern Approach”, 3rd Ed., Prentice Hall, 2010

[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