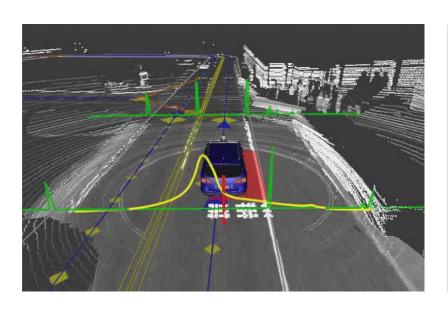
RACIOCÍNIO AUTOMÁTICO PARA DECISÃO SEQUENCIAL

Luís Morgado

ISEL-ADEETC

PROBLEMAS DE DECISÃO SEQUENCIAL

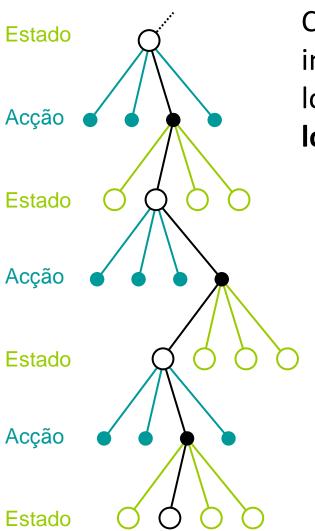




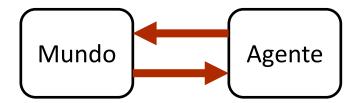




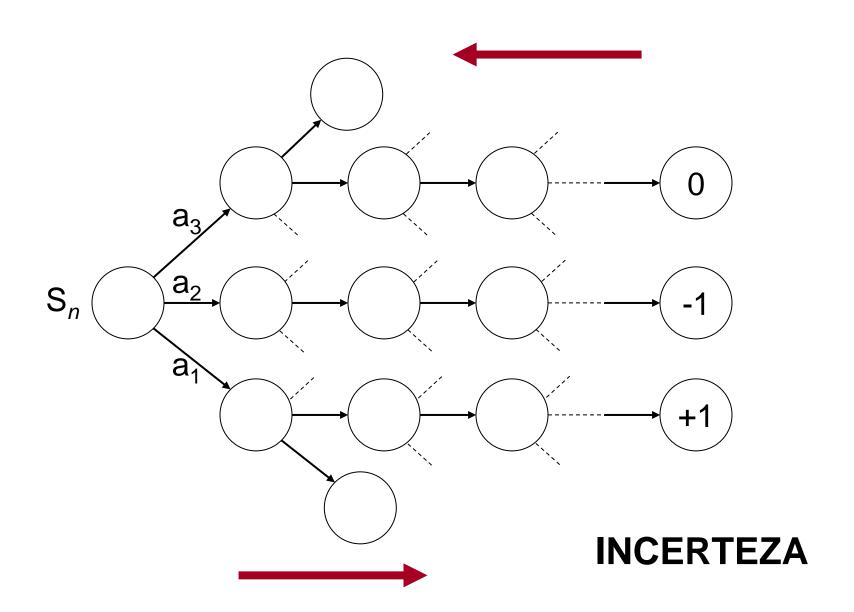
PROCESSOS DE DECISÃO SEQUENCIAL



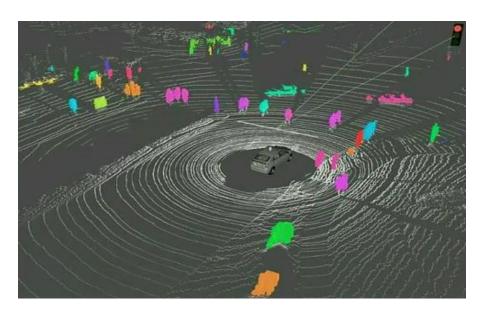
Como prever e controlar o desenrolar da interacção entre agente e ambiente ao longo do tempo para um **objectivo de longo prazo**?



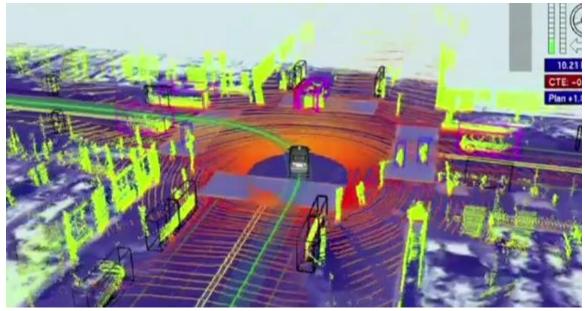
PROCESSOS DE DECISÃO SEQUENCIAL



RACIOCÍNIO COM INCERTEZA

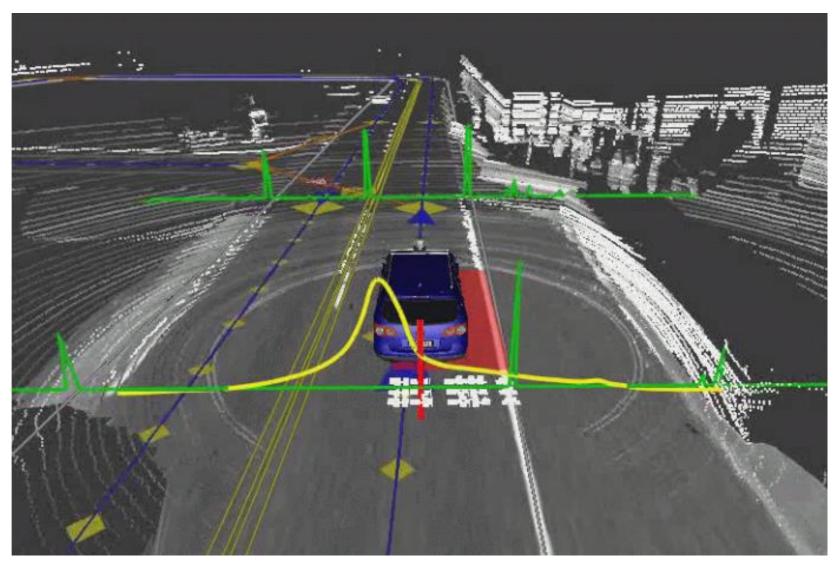








RACIOCÍNIO COM INCERTEZA



PROCESSOS DE DECISÃO SEQUENCIAL

ESPAÇO DE ESTADOS NÃO-DETERMINÍSTICO

Estados

Acções

Transições

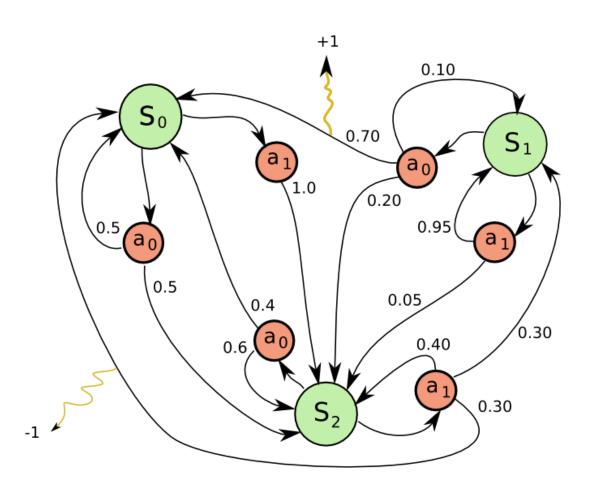
Recompensas

Modelo de Transição

-T(s,a,s')

Modelo de Recompensa

-R(s,a,s')



PROCESSOS DE DECISÃO SEQUENCIAL

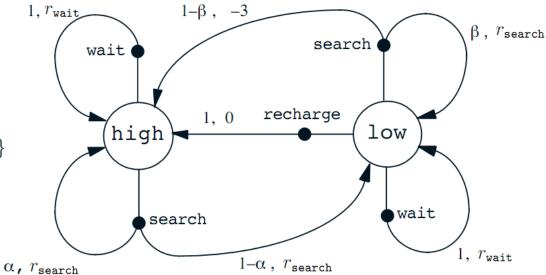
EXEMPLO: ROBOT DE RECICLAGEM

s	s'	a	p(s' s,a)	r(s, a, s')
high	high	search	α	$r_{\mathtt{search}}$
high	low	search	$1-\alpha$	$r_{\mathtt{search}}$
low	high	search	$1-\beta$	-3
low	low	search	β	$r_{\mathtt{search}}$
high	high	wait	1	$r_{\mathtt{wait}}$
high	low	wait	0	$r_{\mathtt{wait}}$
low	high	wait	0	$r_{\mathtt{wait}}$
low	low	wait	1	$r_{\mathtt{wait}}$
low	high	recharge	1	0 1,
low	low	recharge	0	0.
			'	

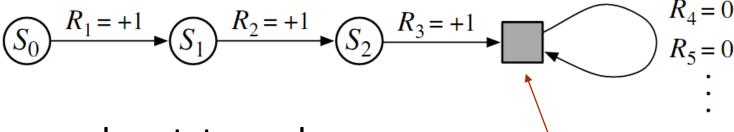
 $\mathbb{S} = \{ \mathtt{high}, \mathtt{low} \}$ (carga da bateria)

 $\mathcal{A}(\texttt{high}) \ = \ \{\texttt{search}, \texttt{wait}\}$

 $A(low) = \{search, wait, recharge\}$



CADEIAS DE MARKOV



Retorno com desconto temporal:

Estado final (absorvente)

$$G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}, \quad \gamma \in [0,1]$$
 - Factor de desconto

PROPRIEDADE DE MARKOV

Um processo estocástico tem a **propriedade de Markov** se a distribuição probabilística condicional dos estados futuros de um processo **depender exclusivamente do estado presente**

RETORNO DE HORIZONTE INFINITO

- Não está limitado a uma gama finita de valores
- Necessário ponderar a distância no tempo das recompensas
 - Recompensas descontadas (discounted reward)
 - Factor de desconto $\gamma \in [0,1]$

$$R_t = r_{t+1} + \gamma \cdot r_{t+2} + \gamma^2 \cdot r_{t+3} + \gamma^3 \cdot r_{t+4} + \dots = \sum_{i=1}^{\infty} \gamma^{i-1} \cdot r_{t+i}$$

Recompensas descontadas no tempo

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

 γ – taxa de desconto para recompensas diferidas no tempo

t = 0, 1, 2, ... - tempo discreto

$$S_{t} \xrightarrow{a_{t}} S_{t+1} \xrightarrow{r_{t+2}} S_{t+2} \xrightarrow{r_{t+3}} S_{t+3} \xrightarrow{a_{t+3}} S_{t+3}$$

Cadeia de Markov

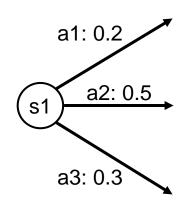
POLÍTICA COMPORTAMENTAL

- Forma de representação do comportamento do agente
- Define qual a acção que deve ser realizada em cada estado
- Política determinista

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

Política não determinista

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



Objectivo

- Maximizar o valor (retorno) de uma sequência de acções
 - Política comportamental
 - Valor de um estado com base numa política

$$V^{\pi}(s) = E\langle r_1 + \gamma r_2 + \gamma^2 r_3 + ... | s_0 = s, \pi \rangle$$

Factor de desconto $\gamma \in [0,1]$ para recompensas diferidas no tempo

- Obter política óptima
- Definir uma relação de ordem parcial entre políticas, existindo, pelo menos, uma política óptima

$$\pi^* = \arg\max_{\pi} V^{\pi}$$

EXEMPLO

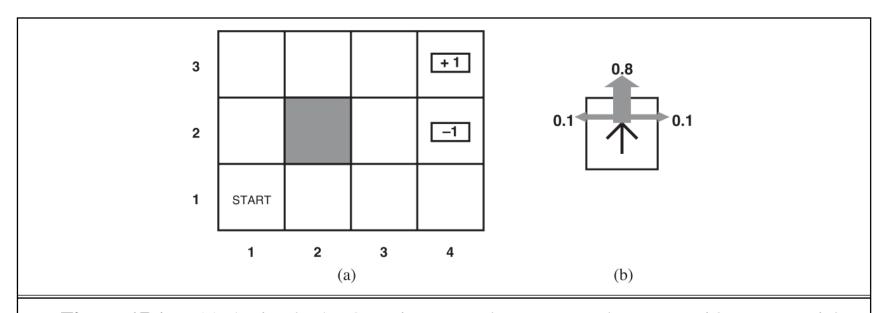


Figure 17.1 (a) A simple 4×3 environment that presents the agent with a sequential decision problem. (b) Illustration of the transition model of the environment: the "intended" outcome occurs with probability 0.8, but with probability 0.2 the agent moves at right angles to the intended direction. A collision with a wall results in no movement. The two terminal states have reward +1 and -1, respectively, and all other states have a reward of -0.04.

EXEMPLO

POLÍTICA COMPORTAMENTAL

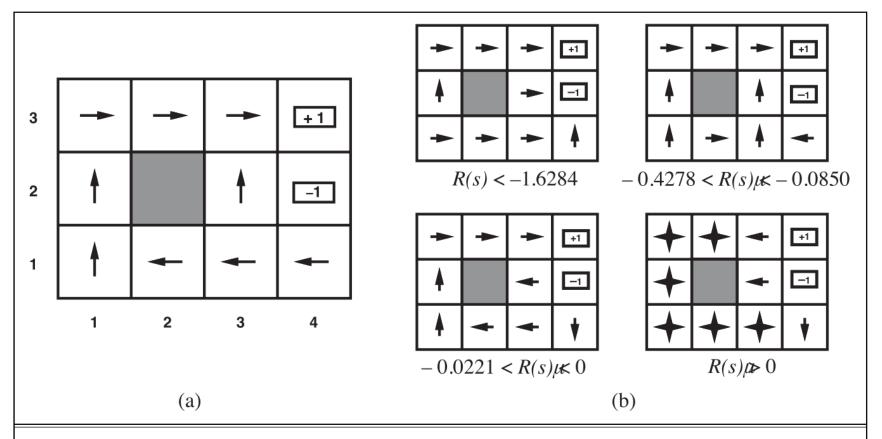


Figure 17.2 (a) An optimal policy for the stochastic environment with R(s) = -0.04 in the nonterminal states. (b) Optimal policies for four different ranges of R(s).

EXEMPLO

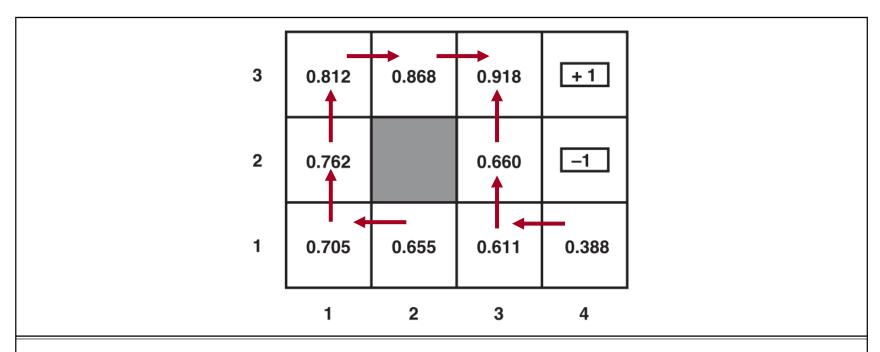


Figure 17.3 The utilities of the states in the 4×3 world, calculated with $\gamma = 1$ and R(s) = -0.04 for nonterminal states.

- Utilidade (valor) de estado
 - Medida de valor de estado
 - Reflecte congruência com uma finalidade definida (objectivo)
 - Depende da **política** comportamental utilizada
- R(s)
 - Recompensa a curto-prazo
- U(s)V(s)
 - Recompensa a longo-prazo

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

$$U^{\pi}(s) = E\langle r_1 + \gamma r_2 + \gamma^2 r_3 + ... \rangle$$
 Equações de Bellman
$$= E\langle r_1 + \gamma U^{\pi}(s') \rangle$$

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

Valor esperado

Cadeia de Markov

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

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

Política:
$$\pi$$
 s
 s'
 s''
 s''
 s''
 s''

Episódio

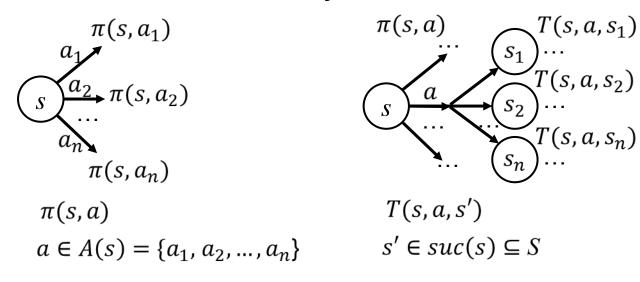
1 r_1^1 r_2^1 r_3^1 ...

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

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

Política

Transição de estado com base num modelo



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

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^*

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

PROGRAMAÇÃO DINÂMICA

Resolução de um problema por composição iterativa de soluções parciais Iterar sucessivamente sobre todos os estados, actualizando cada estado com uma aproximação do valor óptimo

Iteração do Valor 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:

Política de selecção de acção greedy

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

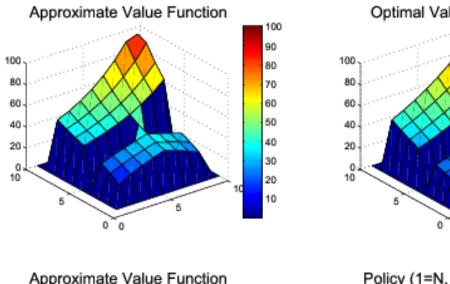
PROGRAMAÇÃO DINÂMICA

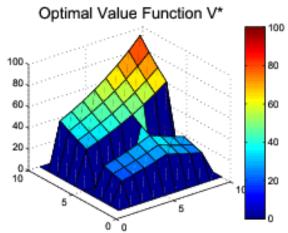
- Iteração
 - Backup
 - Utilizar retornos futuros para estimar o retorno do estado actual

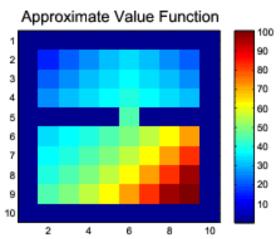
$$V(s) := \max_{a} \sum_{s'} T_{ss'}^{a} \left[R_{ss'}^{a} + \gamma V(s') \right] \forall s \in S$$

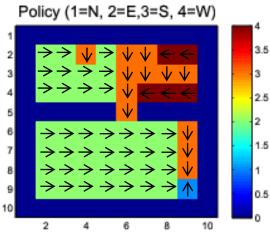
$$T_{ss'}^{a} \qquad T_{ss''}^{a}$$

$$S \stackrel{?}{\sim} R_{ss'}^{a} \qquad R_{ss''}^{a} \qquad S \stackrel{?}{\sim} S$$









CÁLCULO DA UTILIDADE DE ESTADO

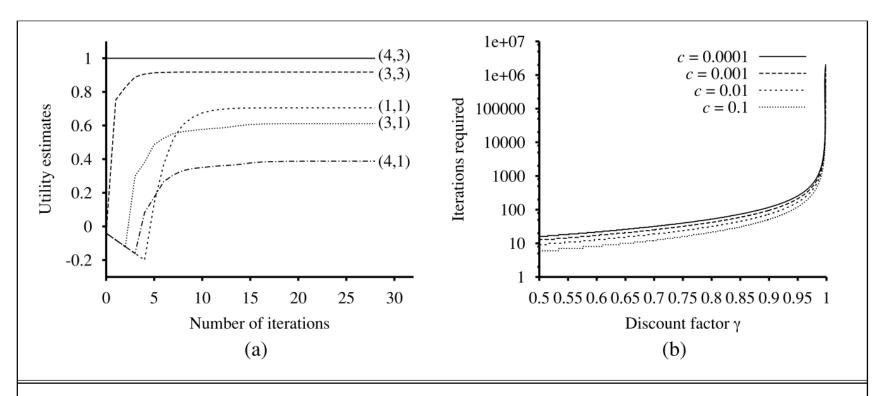


Figure 17.5 (a) Graph showing the evolution of the utilities of selected states using value iteration. (b) The number of value iterations k required to guarantee an error of at most $\epsilon = c \cdot R_{\text{max}}$, for different values of c, as a function of the discount factor γ .

CÁLCULO DA UTILIDADE DE ESTADO

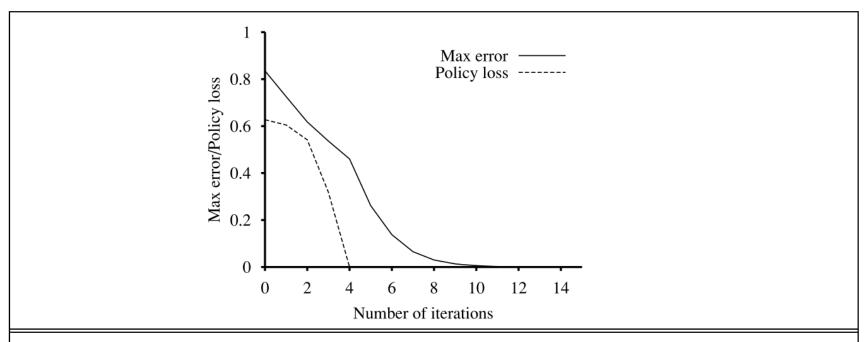


Figure 17.6 The maximum error $||U_i - U||$ of the utility estimates and the policy loss $||U^{\pi_i} - U||$, as a function of the number of iterations of value iteration.

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

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

$$\Delta = \lVert U_{i+1} - U_i
Vert, \quad \Delta < \Delta_{ ext{max}}$$
 (limiar de convergência)

ALGORITMO DE ITERAÇÃO DE VALOR

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation Initialize V(s), for all $s \in S^+$, arbitrarily except that V(terminal) = 0

```
Loop:
```

Loop for each
$$s \in S$$
:
 $v \leftarrow V(s)$

$$V(s) \leftarrow \max_{a} \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until
$$\Delta < \theta$$

 $\Delta \leftarrow 0$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \operatorname{arg\,max}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

[Sutton & Barto, 2020]

$$T(s, a, s') = p(s'|s, a)$$

 $R(s, a, s') = p(r|s, a)$

UTILIDADE DE ESTADO

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

 Iteração de valor pode ser vista como uma propagação de informação de valor através do espaço de estados

Síncrona

- As iterações actualizam sistematicamente cada estado
- Manutenção de uma cópia da utilidade anterior

Assíncrona

- As iterações não actualizam sistematicamente cada estado
- Necessário continuar a iterar até que todos os estados tenham sido actualizados

- 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
 - \bullet 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

Problemas

- Dimensão dos espaços de estados
 - Complexidade computacional exponencial
- Dificuldade de definição das dinâmicas (por exemplo a partir de dados experimentais)
- Dinâmicas desconhecidas

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

REFERÊNCIAS

[Russel & Norvig, 2010]

S. Russell, P. Norvig, "Artificial Intelligence: A Modern Approach", 3dd Ed., Prentice Hall, 2010

[Sutton & Barto, 1998]

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

[Sutton & Barto, 2012]

R. Sutton, A. Barto, "Reinforcement Learning: An Introduction", 2nd Edition - Preview, MIT Press, 2012

[Sutton & Barto, 2020]

R. Sutton, A. Barto, "Reinforcement Learning: An Introduction", 2nd Edition, MIT Press, 2020

[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