

Aprendizagem por Reforço com Deep Learning, PyTorch e Python

Conteúdo do Curso

O que você aprenderá neste curso:

- Parte 1 – Fundamentos de Aprendizagem por Reforço (Q-Learning)
- Parte 2 – Deep Q-Learning
- Parte 3 – Implementação Deep Q-Learning

Pré-requisitos

- Lógica de programação
 - Programação básica em Python
 - Instalação de softwares
 - Redes Neurais Artificiais (desejável)
 - PyTorch (desejável)
-
- Nível do curso: intermediário
 - Dica: aumentar a velocidade do player
 - Avaliação do curso!

Conteúdo

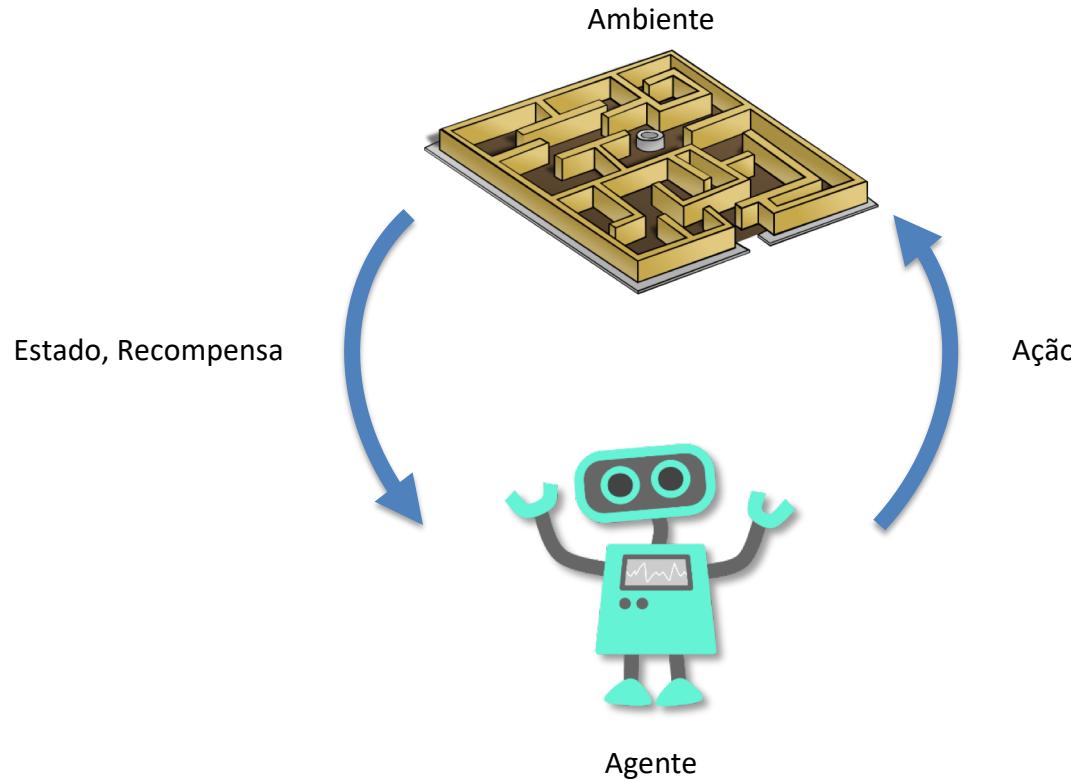
Conteúdo

O que você aprenderá nesta seção:

- O que é aprendizagem por reforço?
- A Equação de Bellman
- O "Plano"
- Markov Decision Process (MDP)
- Política x Plano
- Adição de penalidades ("Living Penalty")
- Q-Learning – Intuição
- Diferença Temporal
- Q-Learning – Visualização

O que é Aprendizagem por Reforço?

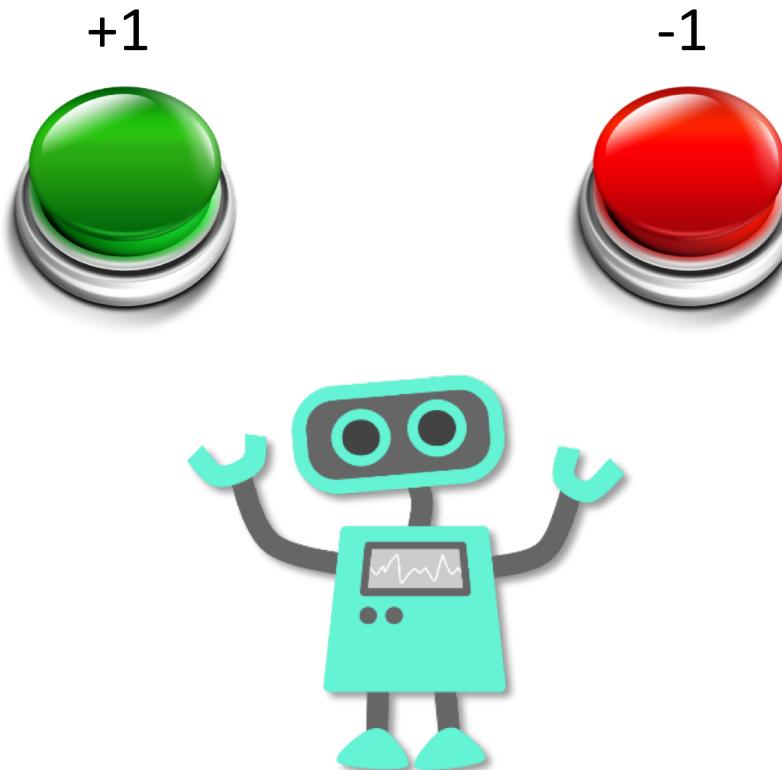
O que é Aprendizagem por Reforço?



O que é Aprendizagem por Reforço?



O que é Aprendizagem por Reforço?



O que é Aprendizagem por Reforço?

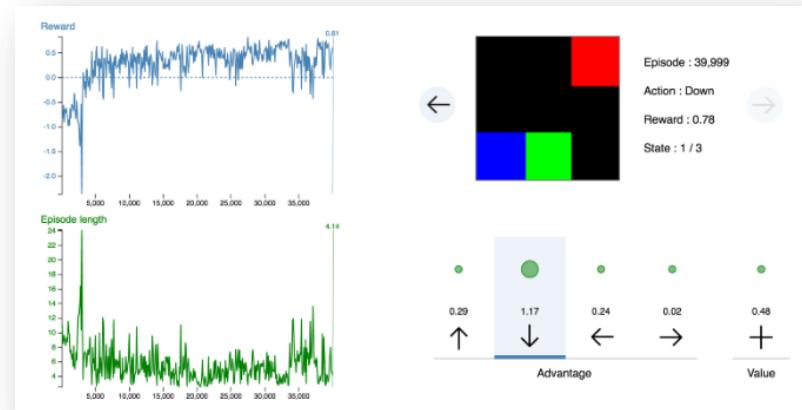


Leitura Adicional

Leitura Adicional:

Simple Reinforcement Learning with Tensorflow (10 Parts)

Arthur Juliani (2016)



Link:

<https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-0-q-learning-with-tables-and-neural-networks-d195264329d0>

Leitura Adicional

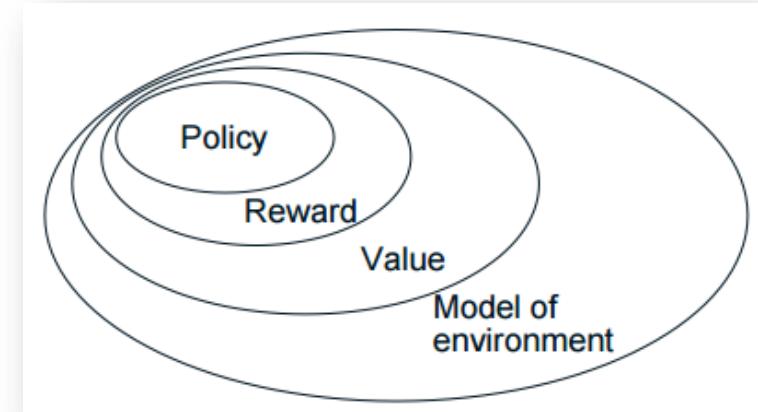
Leitura Adicional:

Reinforcement Learning I: Introduction

Richard Sutton et al. (1998)

Link:

<http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.32.7692>



A Equação de Bellman

A Equação de Bellman

Conceitos:

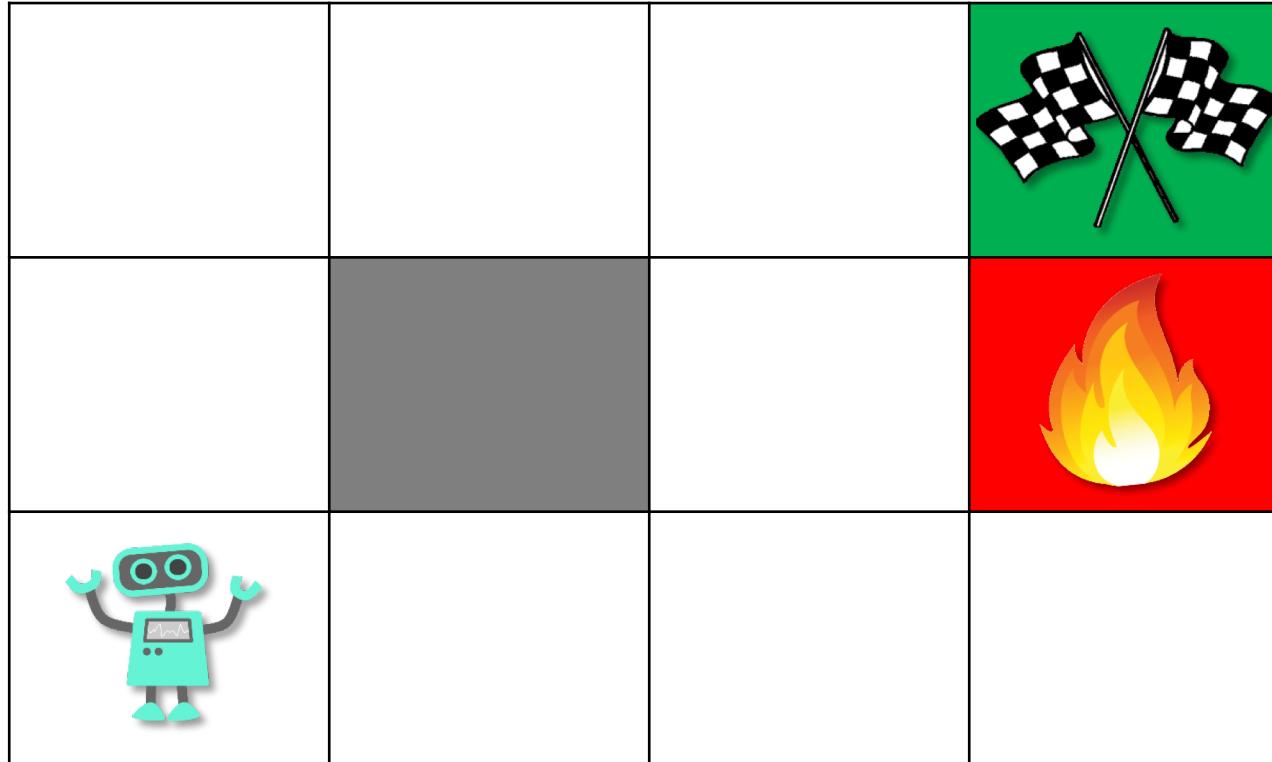
- s – Estado
- a – Ação
- R – Recompensa
- γ – Desconto

A Equação de Bellman

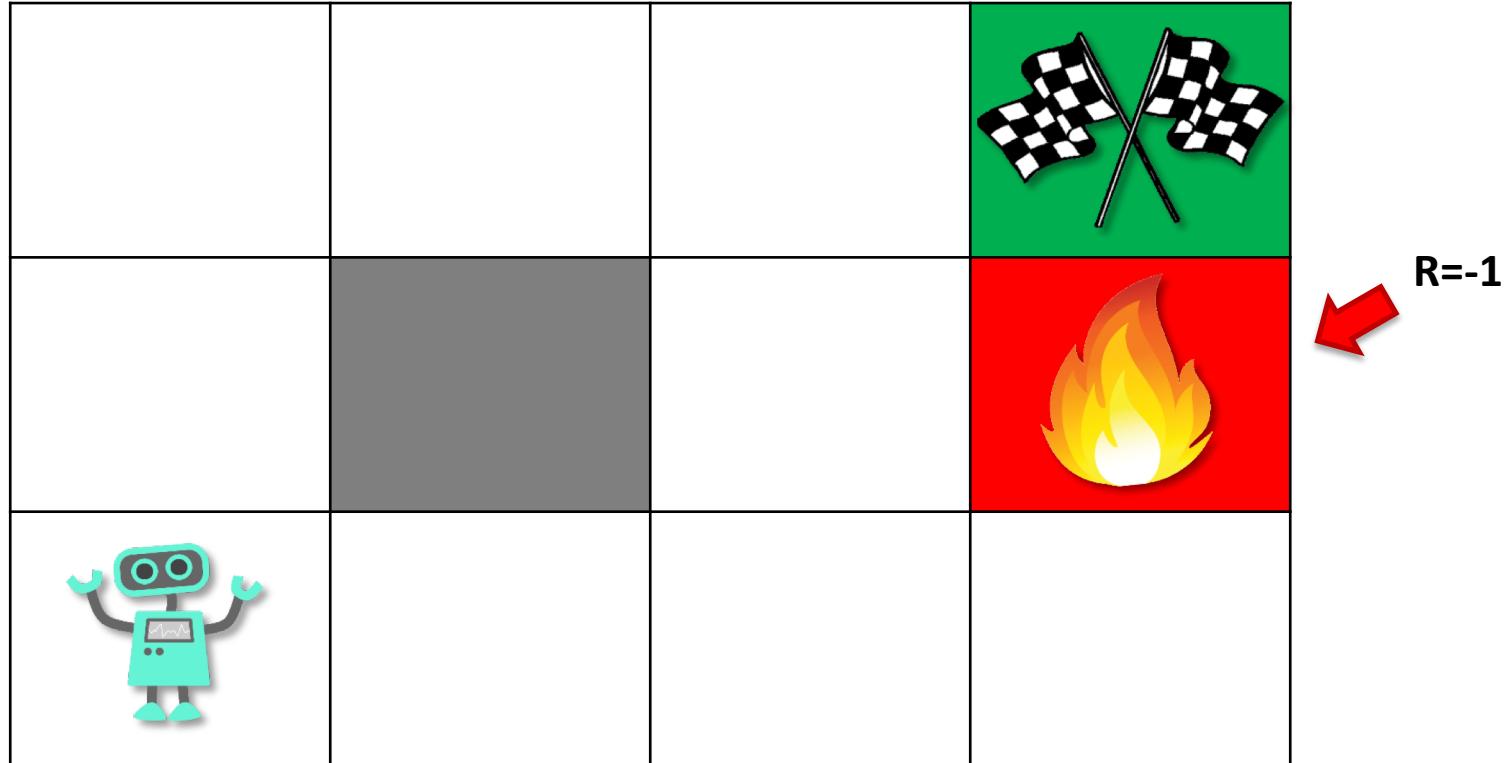


Richard Ernest Bellman

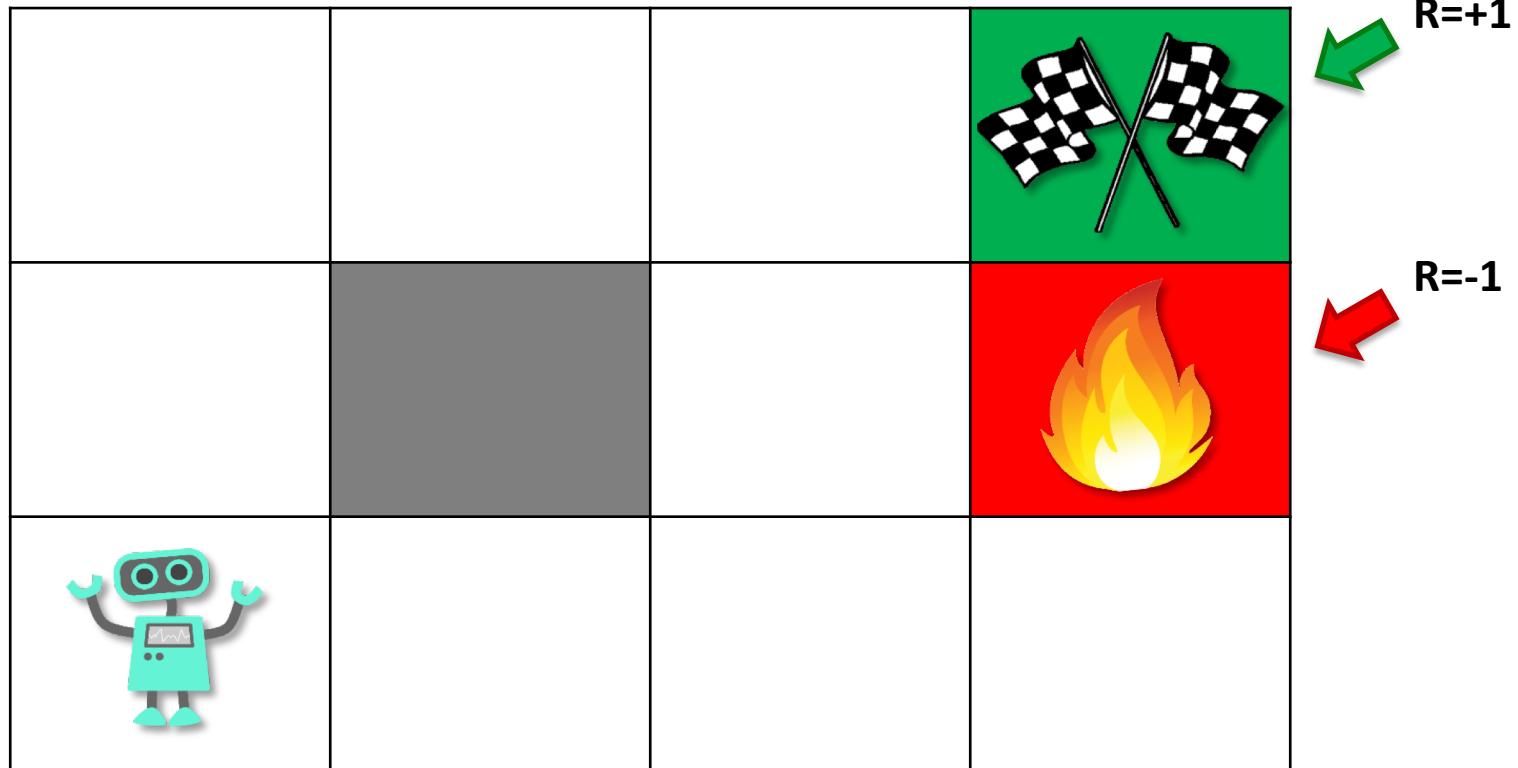
A Equação de Bellman



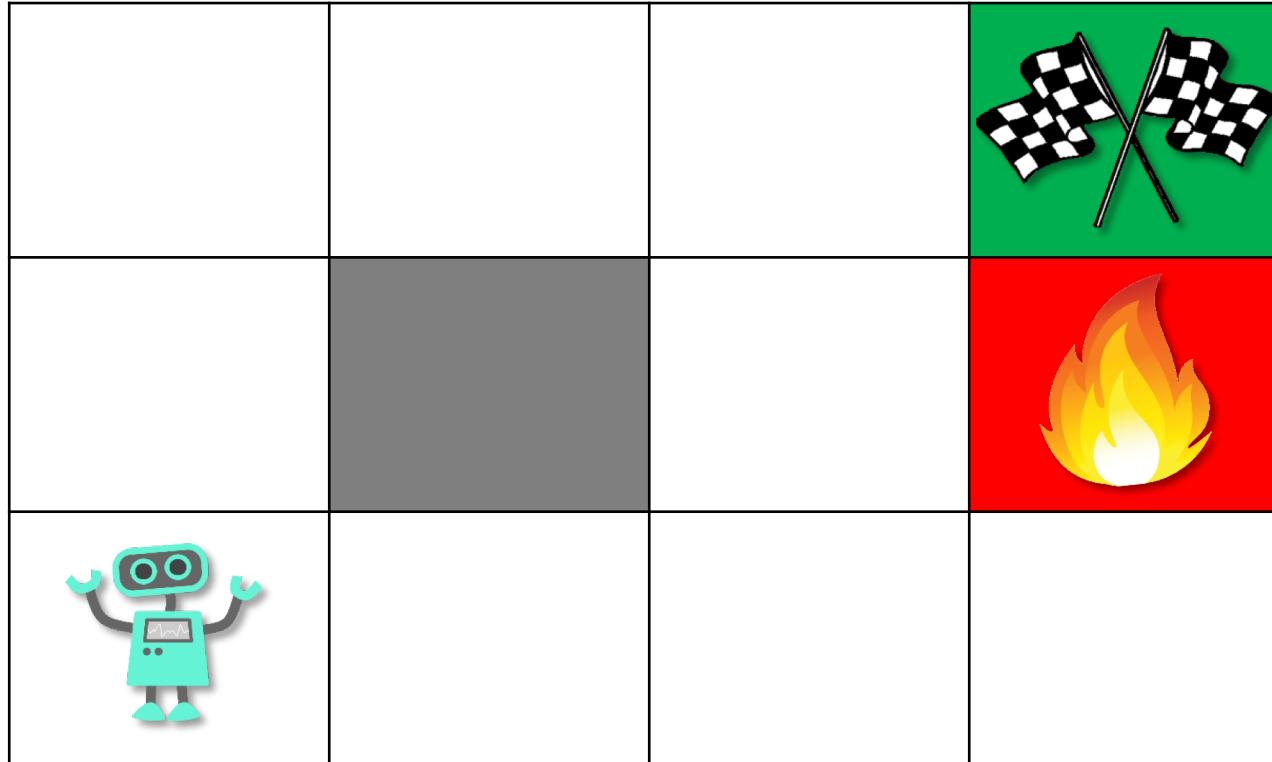
A Equação de Bellman



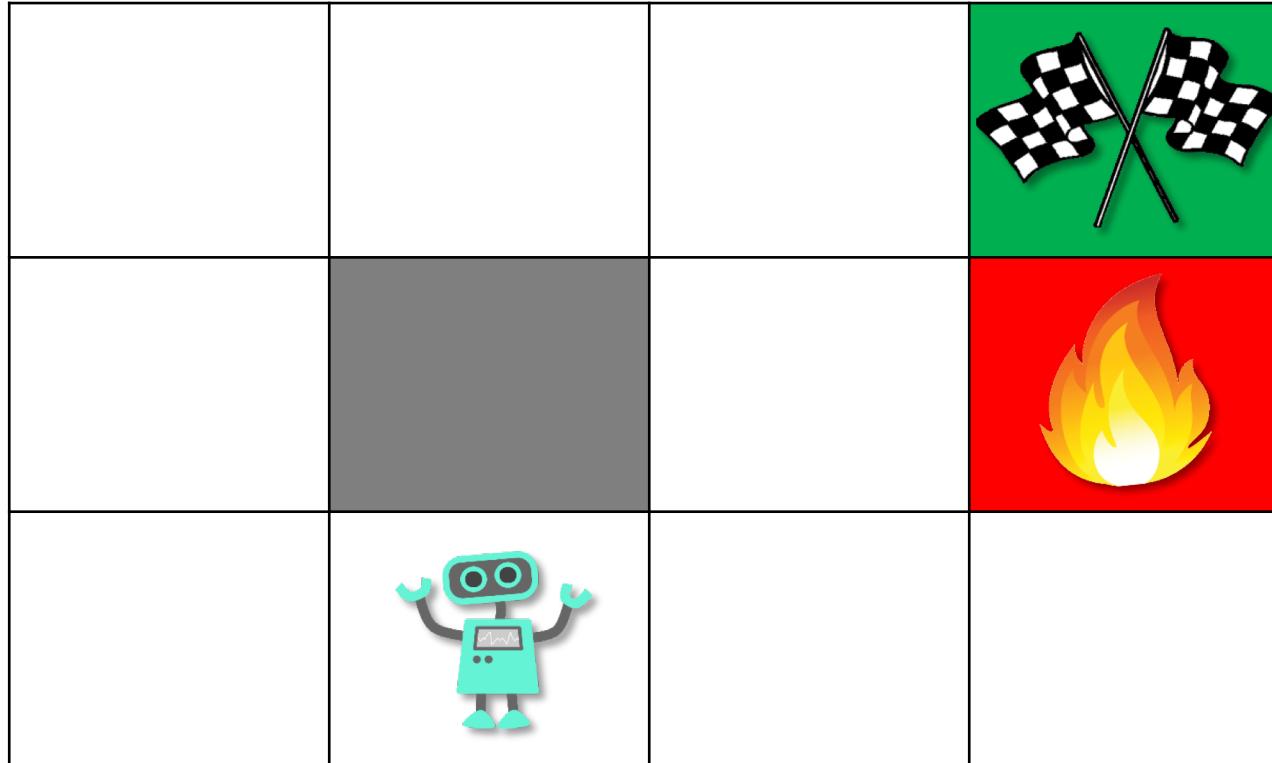
A Equação de Bellman



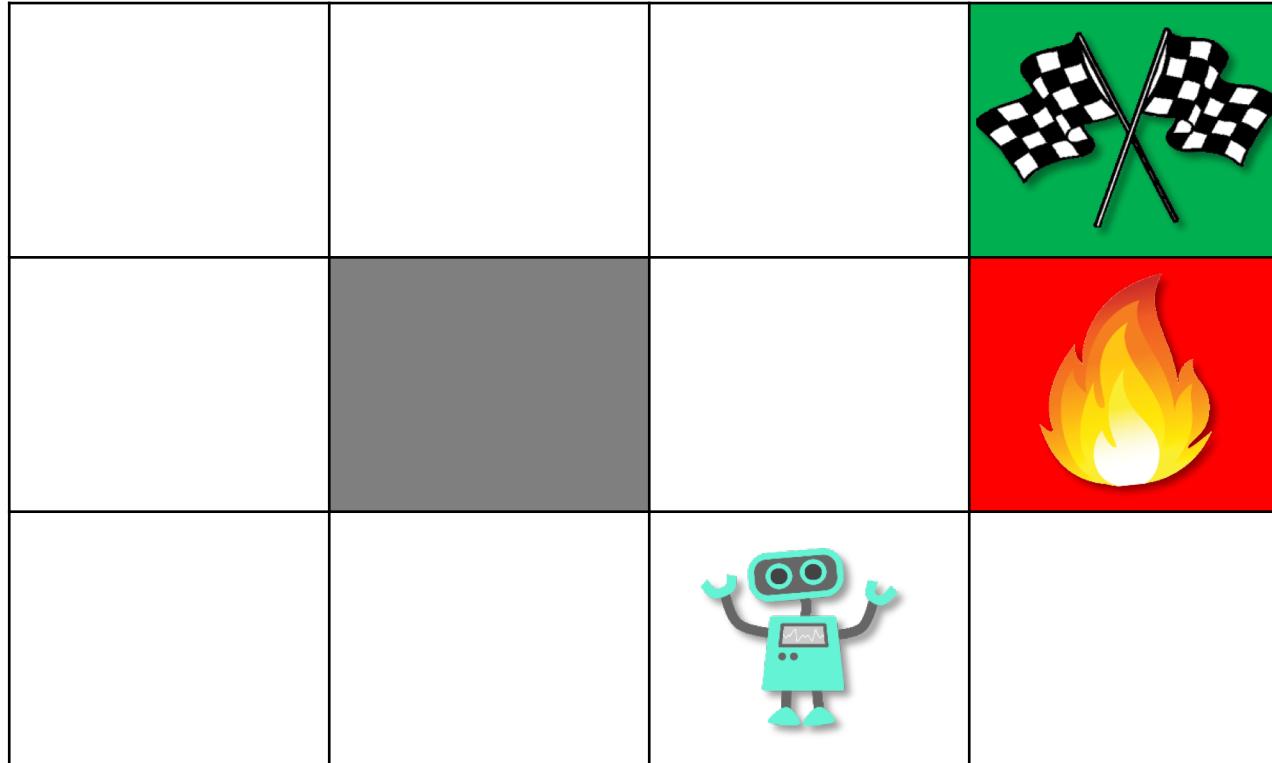
A Equação de Bellman



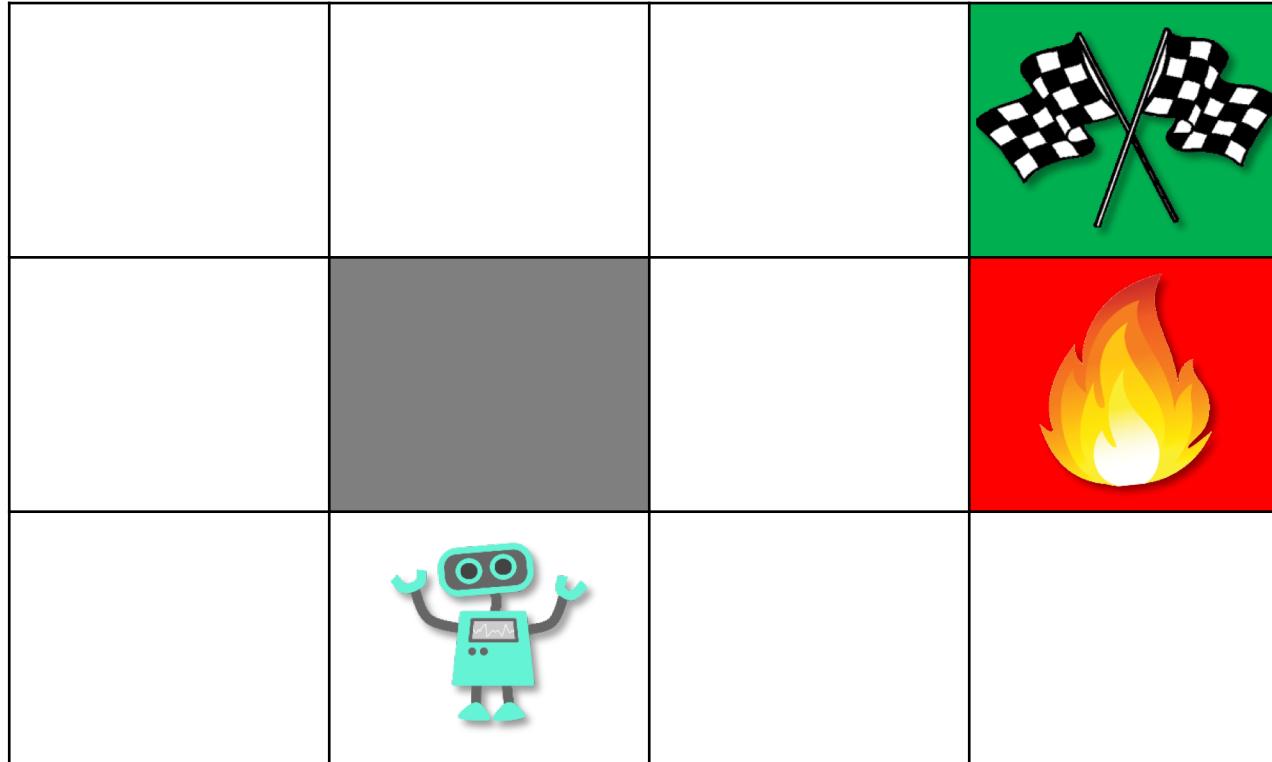
A Equação de Bellman



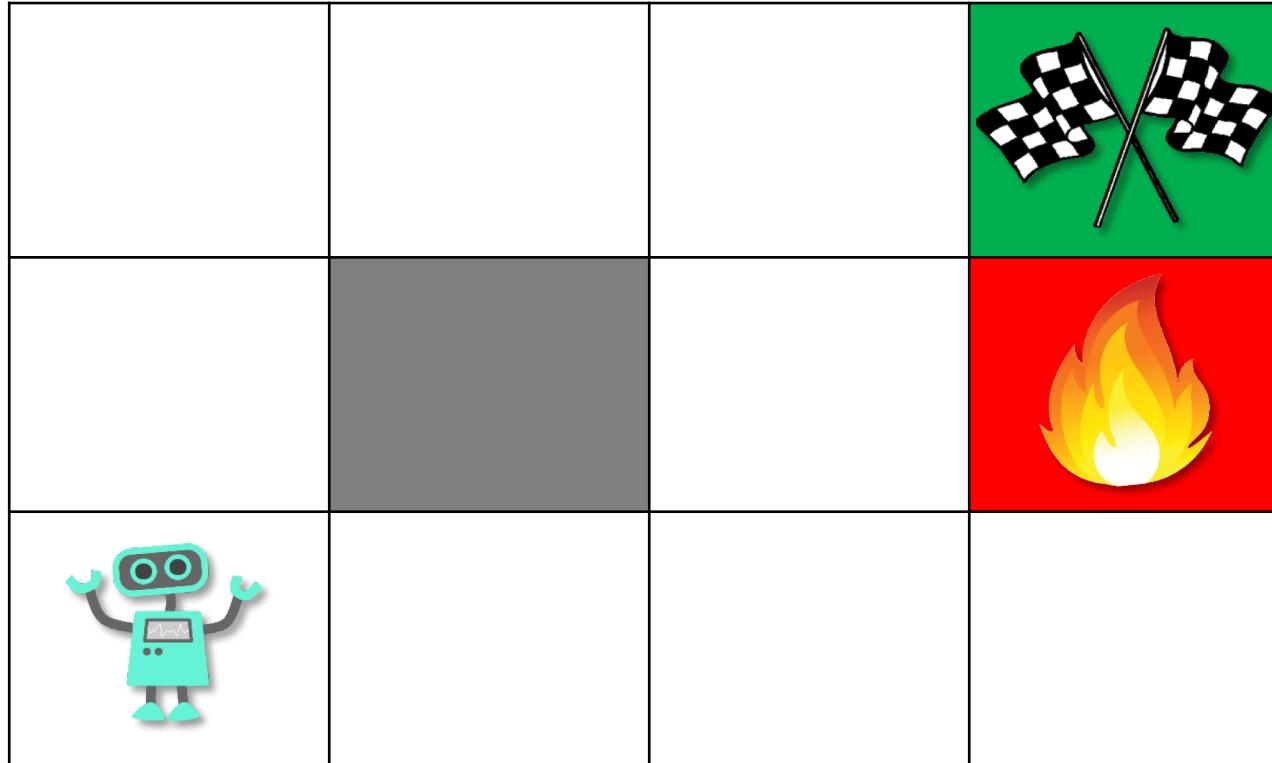
A Equação de Bellman



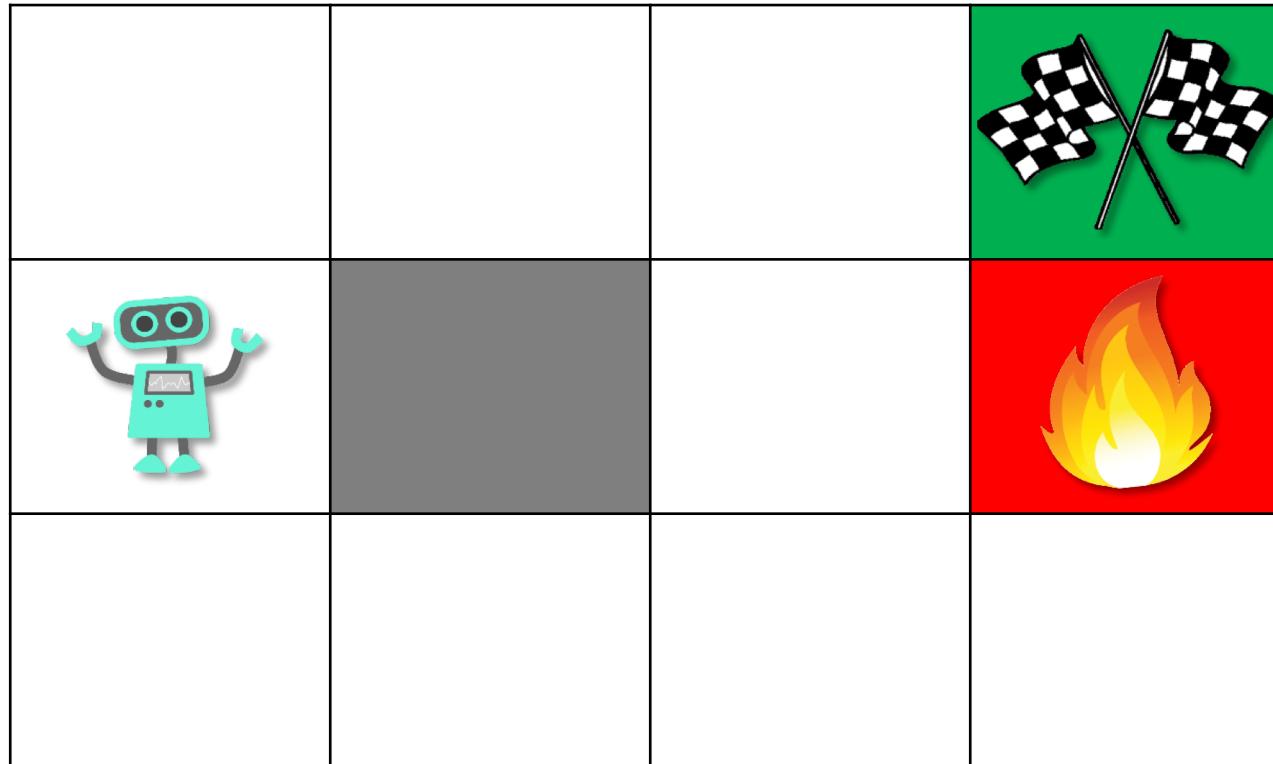
A Equação de Bellman



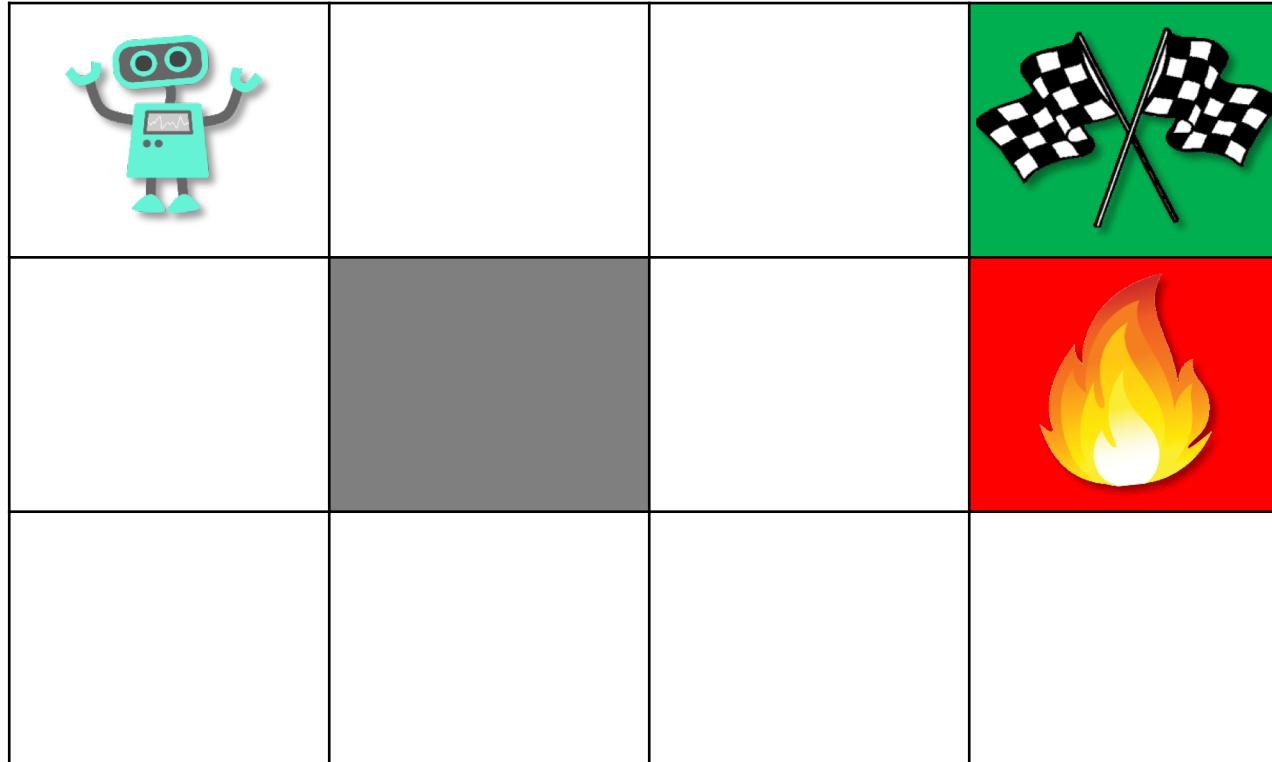
A Equação de Bellman



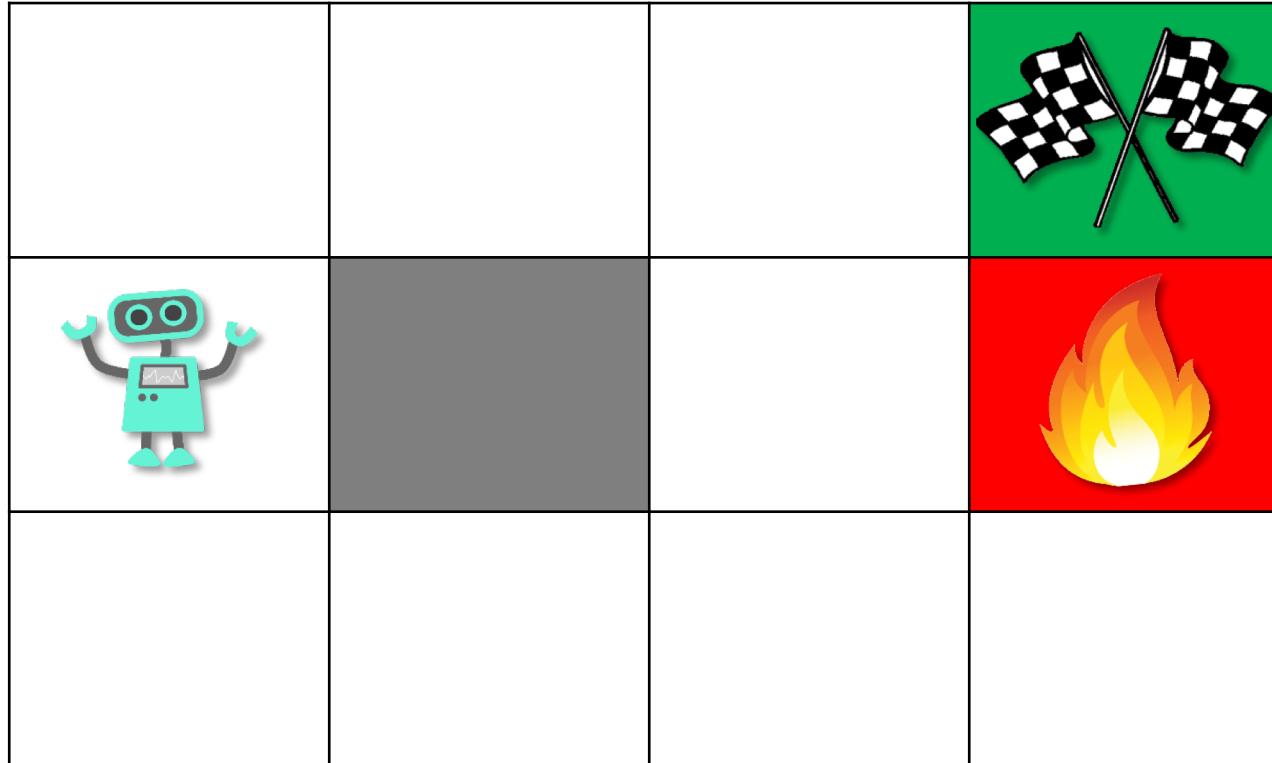
A Equação de Bellman



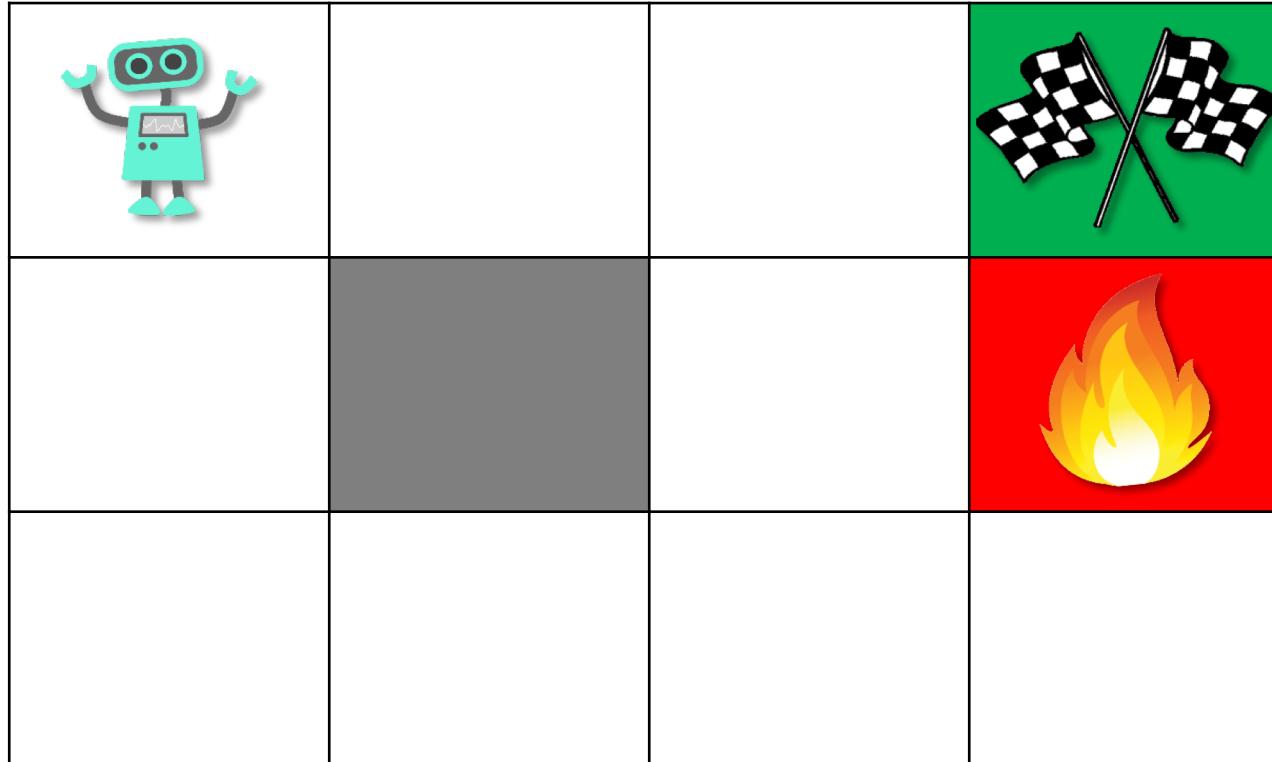
A Equação de Bellman



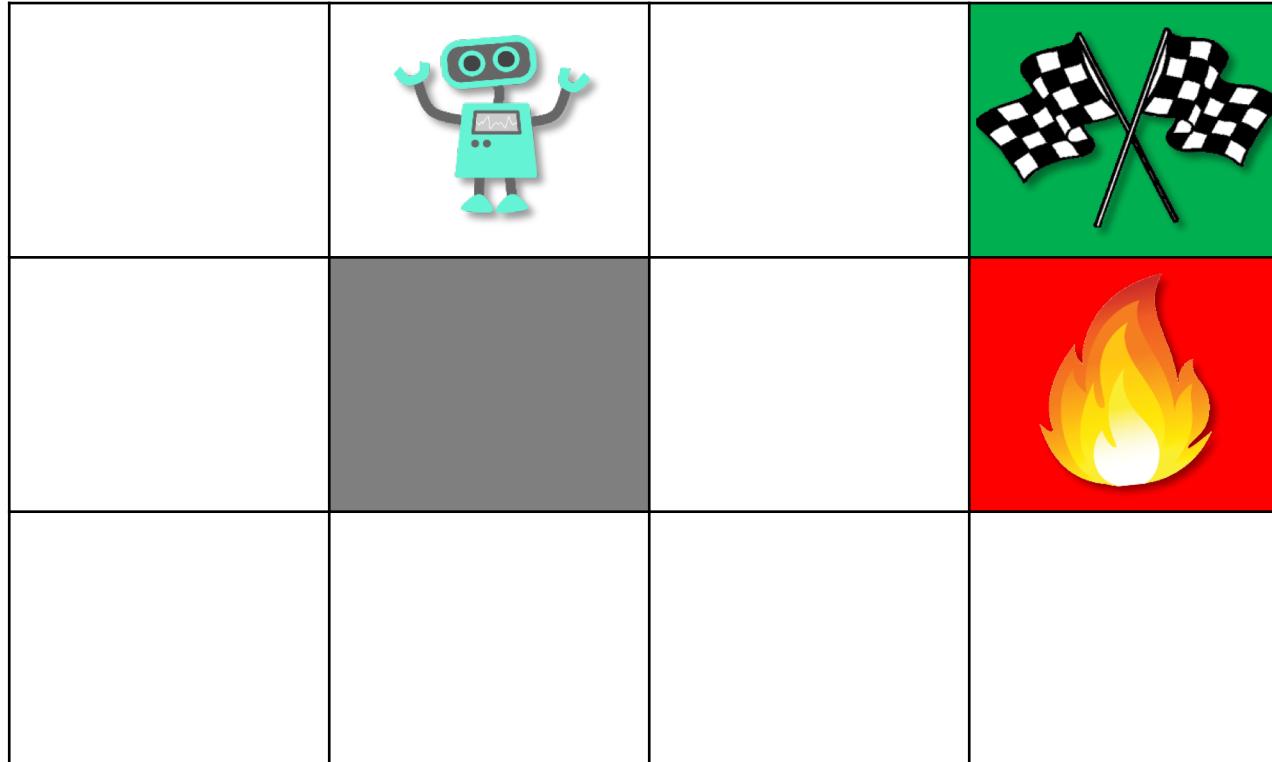
A Equação de Bellman



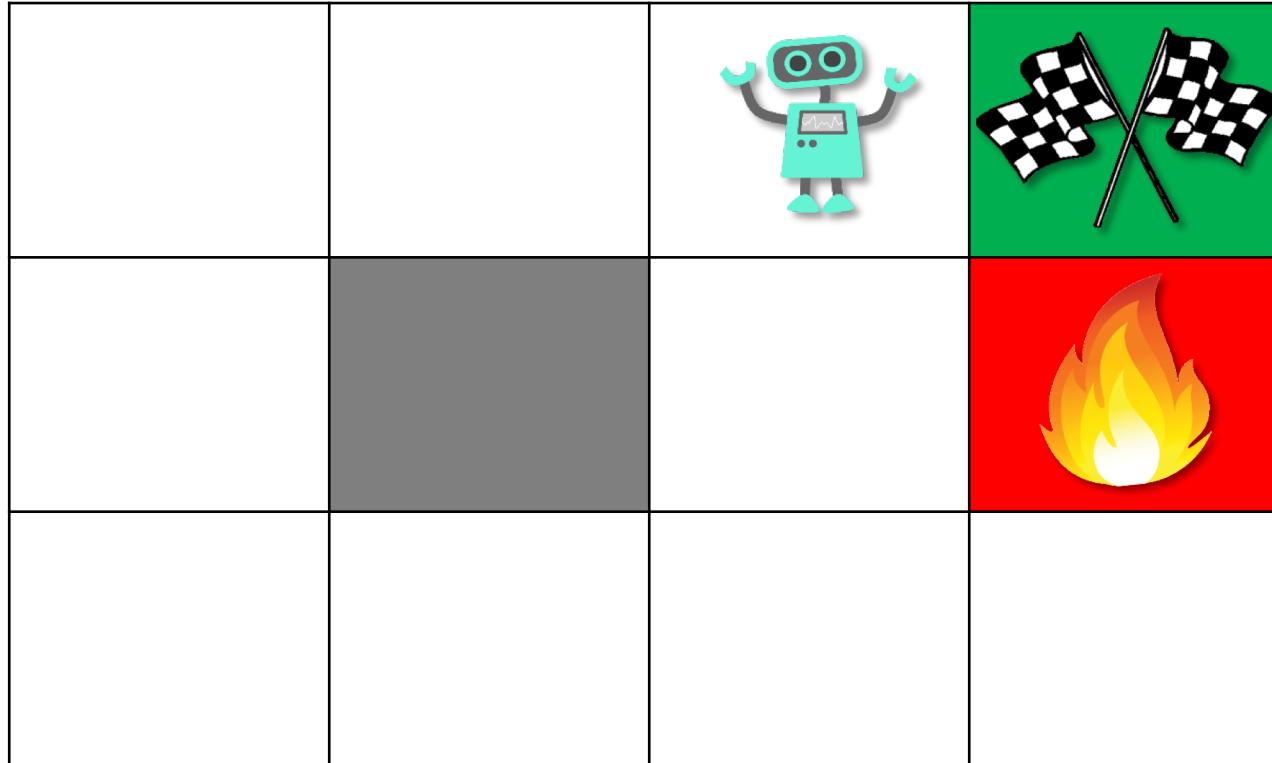
A Equação de Bellman



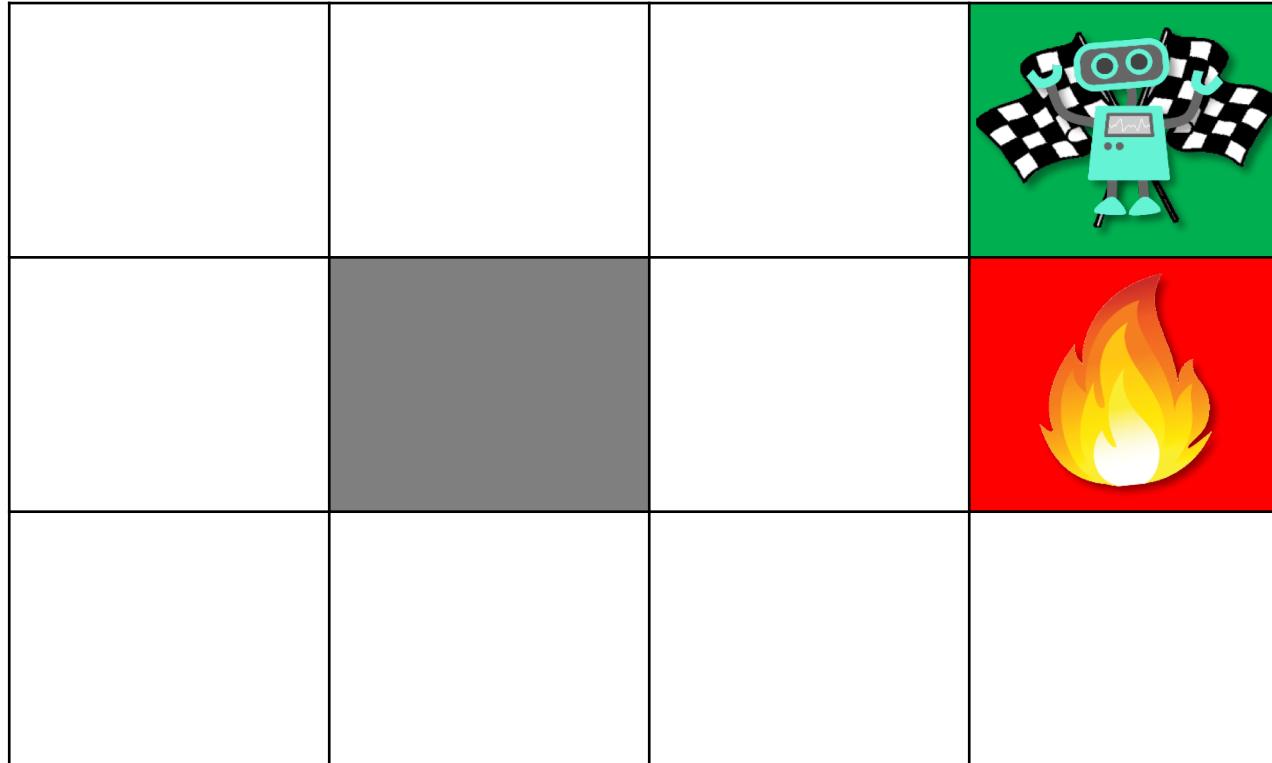
A Equação de Bellman



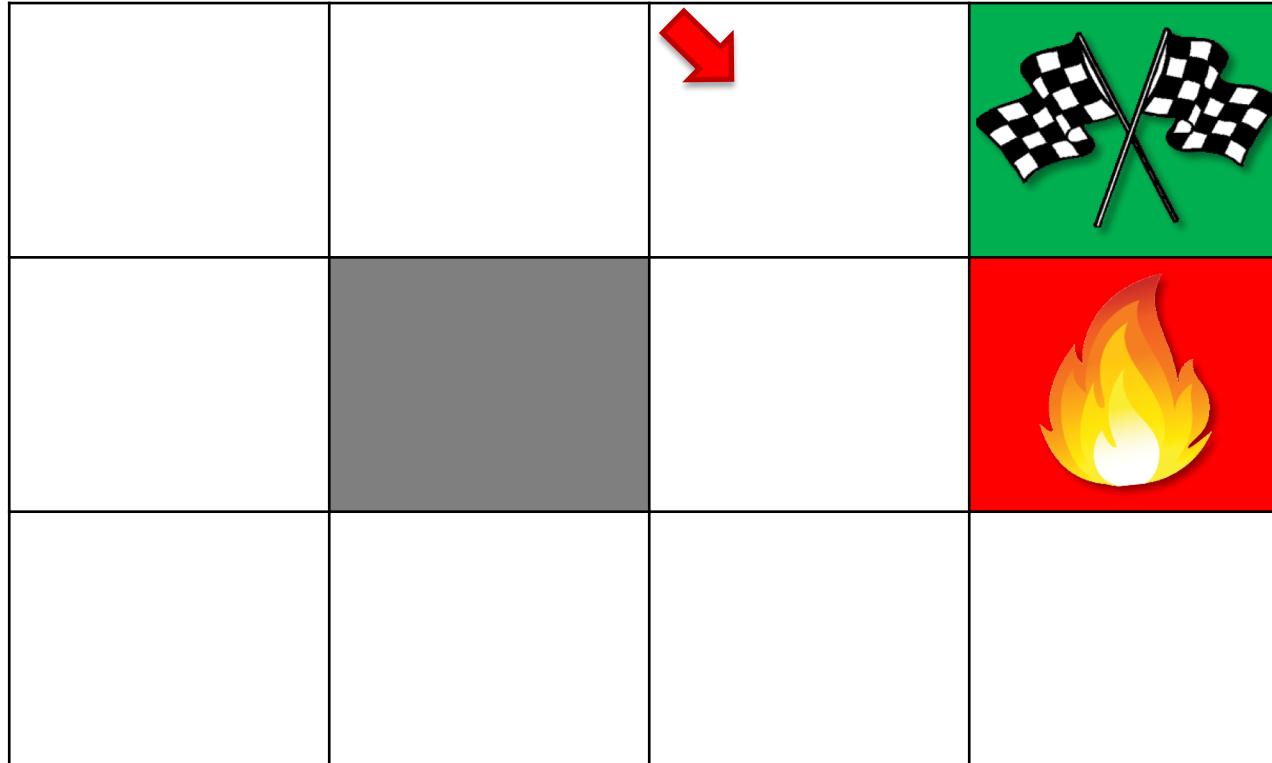
A Equação de Bellman



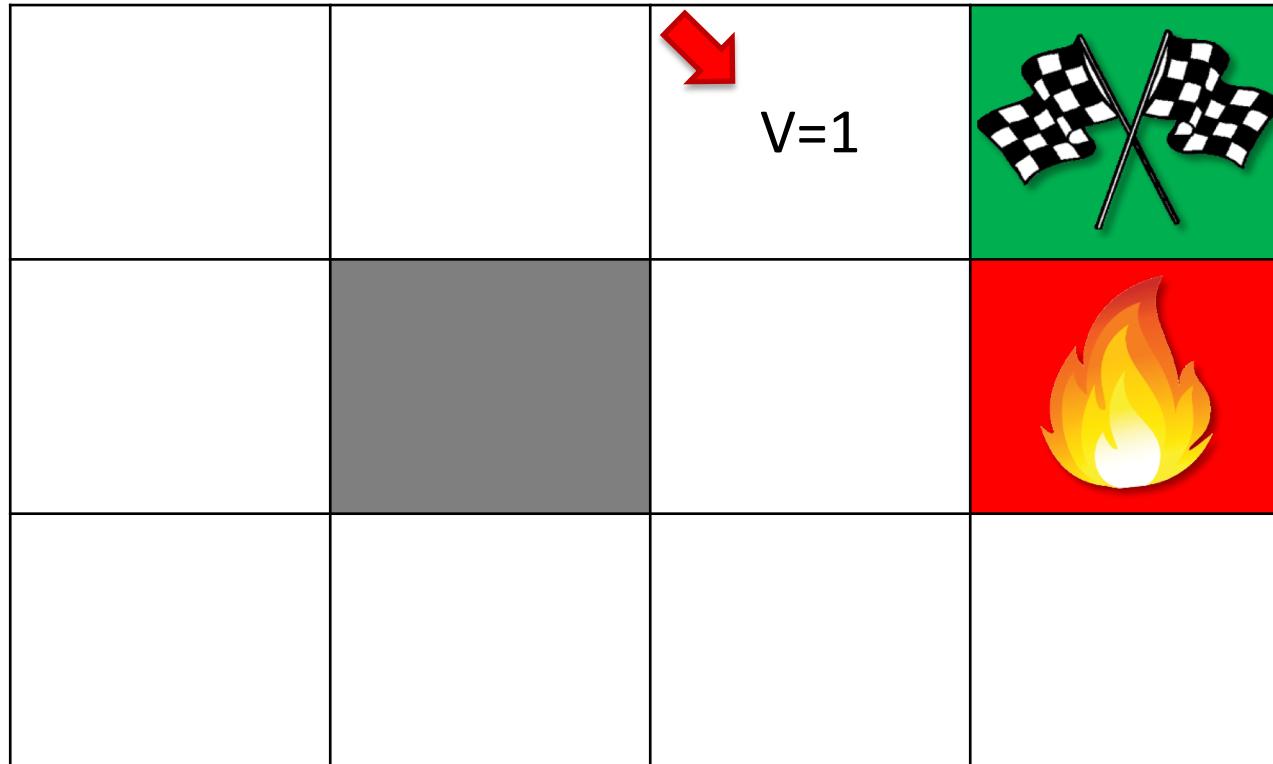
A Equação de Bellman



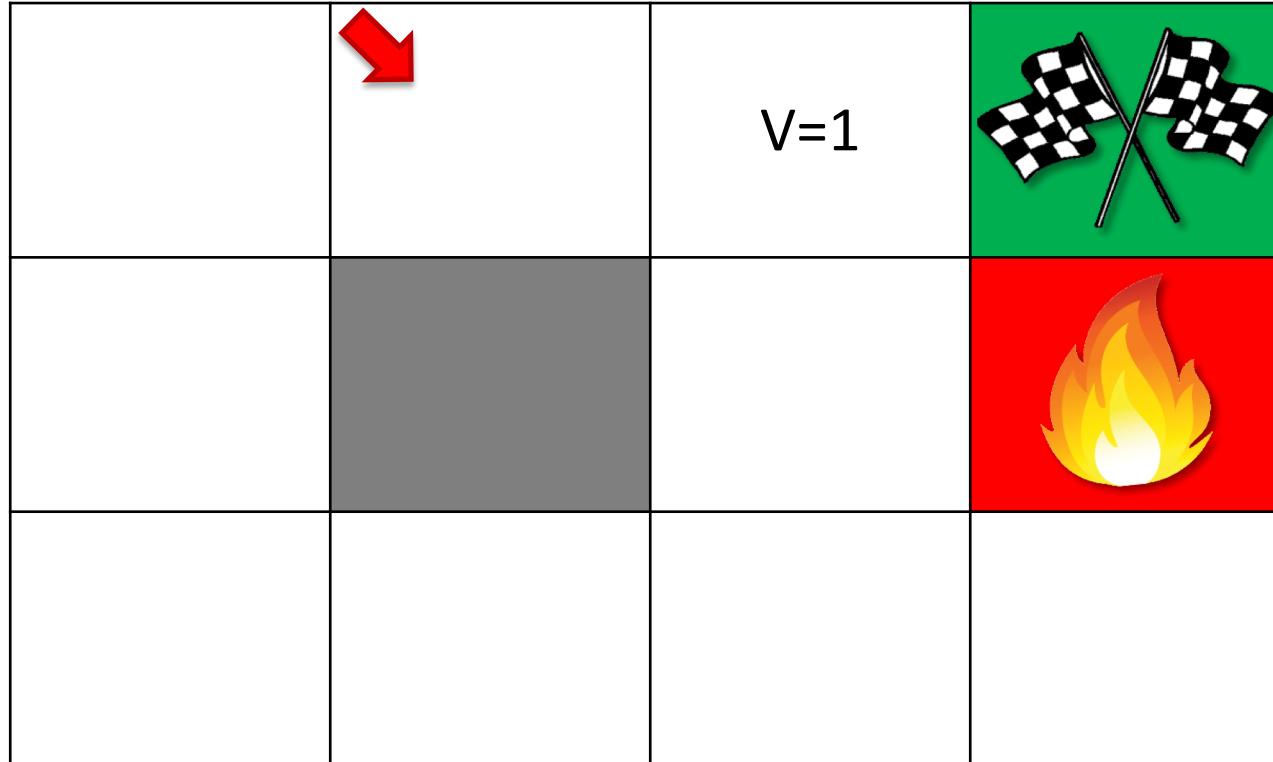
A Equação de Bellman



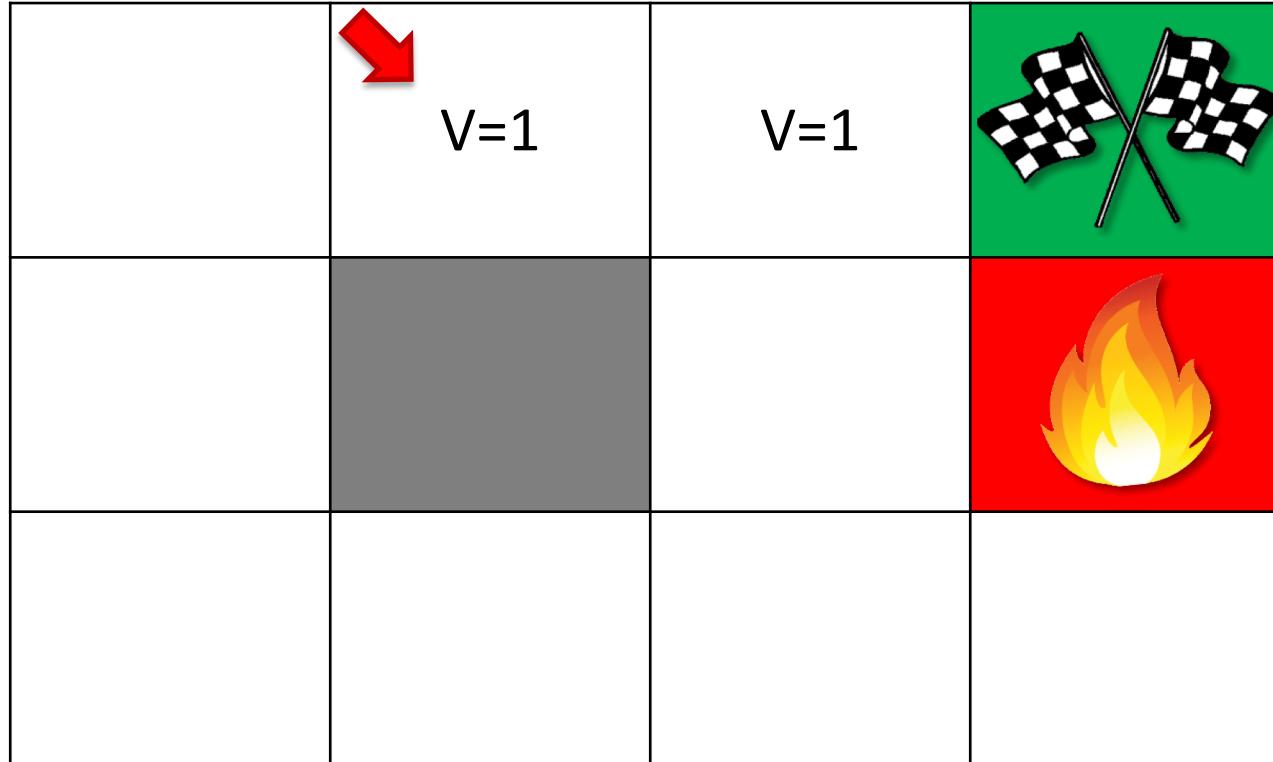
A Equação de Bellman



A Equação de Bellman



A Equação de Bellman



A Equação de Bellman



A Equação de Bellman



A Equação de Bellman



A Equação de Bellman

$V=1$	$V=1$	$V=1$	
 $V=1$			

A Equação de Bellman



A Equação de Bellman

$V=1$	$V=1$	$V=1$	
$V=1$			
 $V=1$			

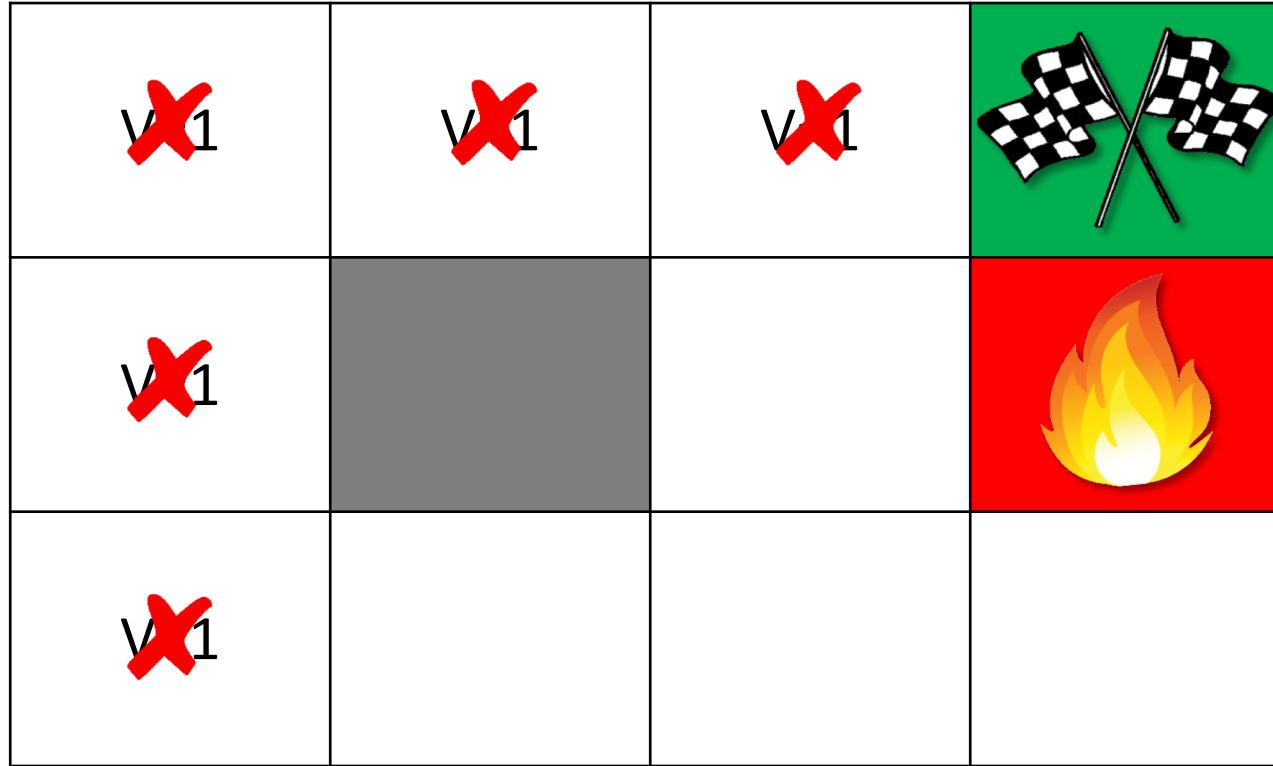
A Equação de Bellman

$V=1$	$V=1$	$V=1$	
$V=1$			
$V=1$			

A Equação de Bellman



A Equação de Bellman



A Equação de Bellman

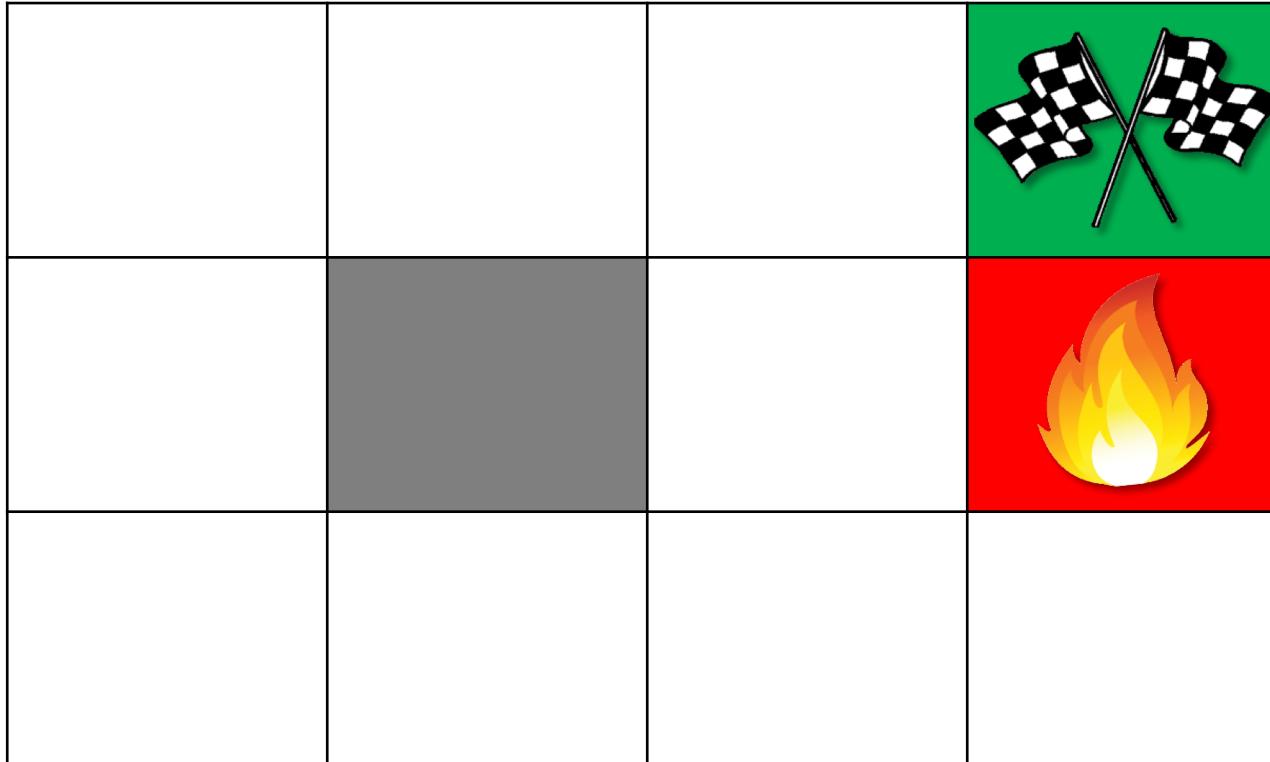


A Equação de Bellman

$$V(s) = \max_a(R(s, a) + \gamma V(s'))$$

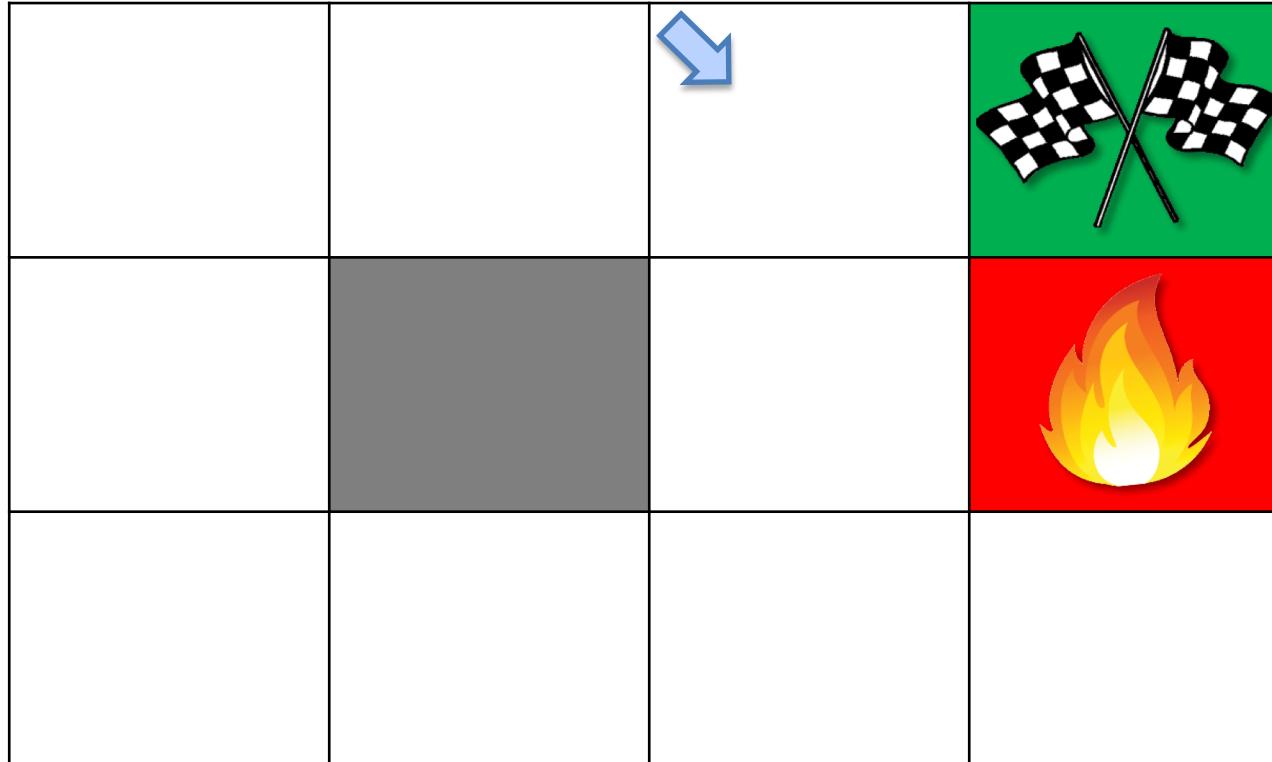
A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



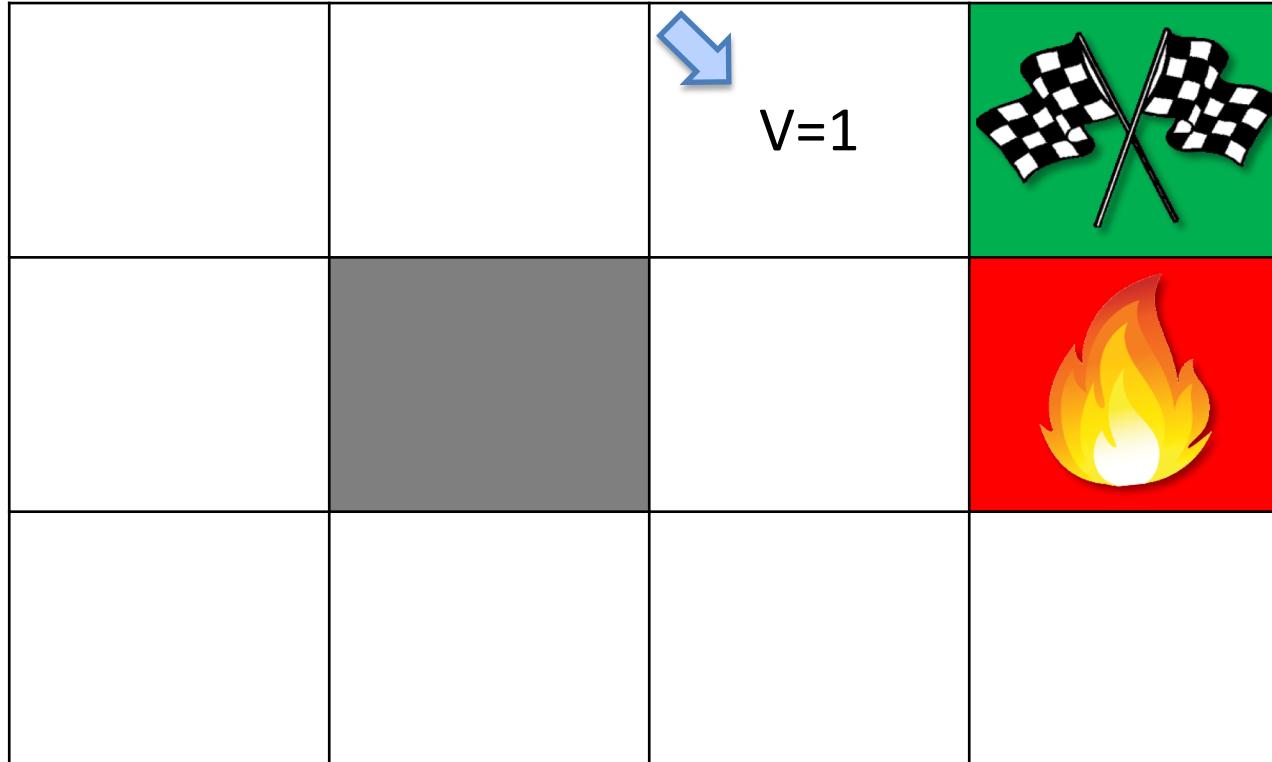
A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



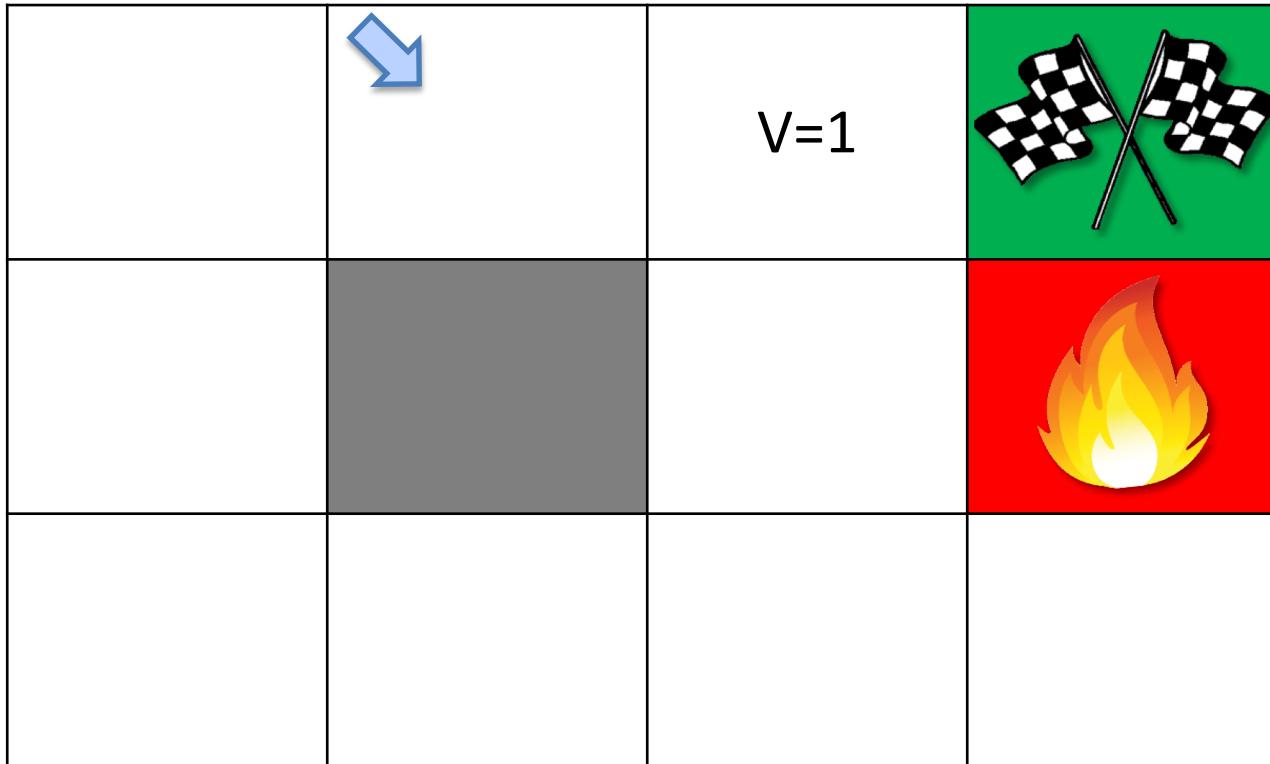
A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

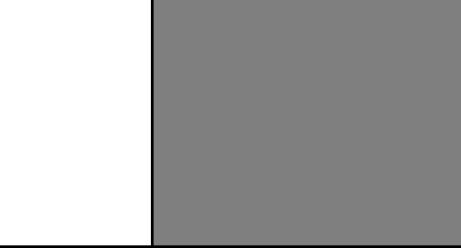
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

	 V=0.9	V=1	
			

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

	V=0.9	V=1	 

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

 V=0.81	V=0.9	V=1	
			

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
			

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
 V=0.73			

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73			
 V=0.66			

$$\gamma = 0.9$$

A Equação de Bellman

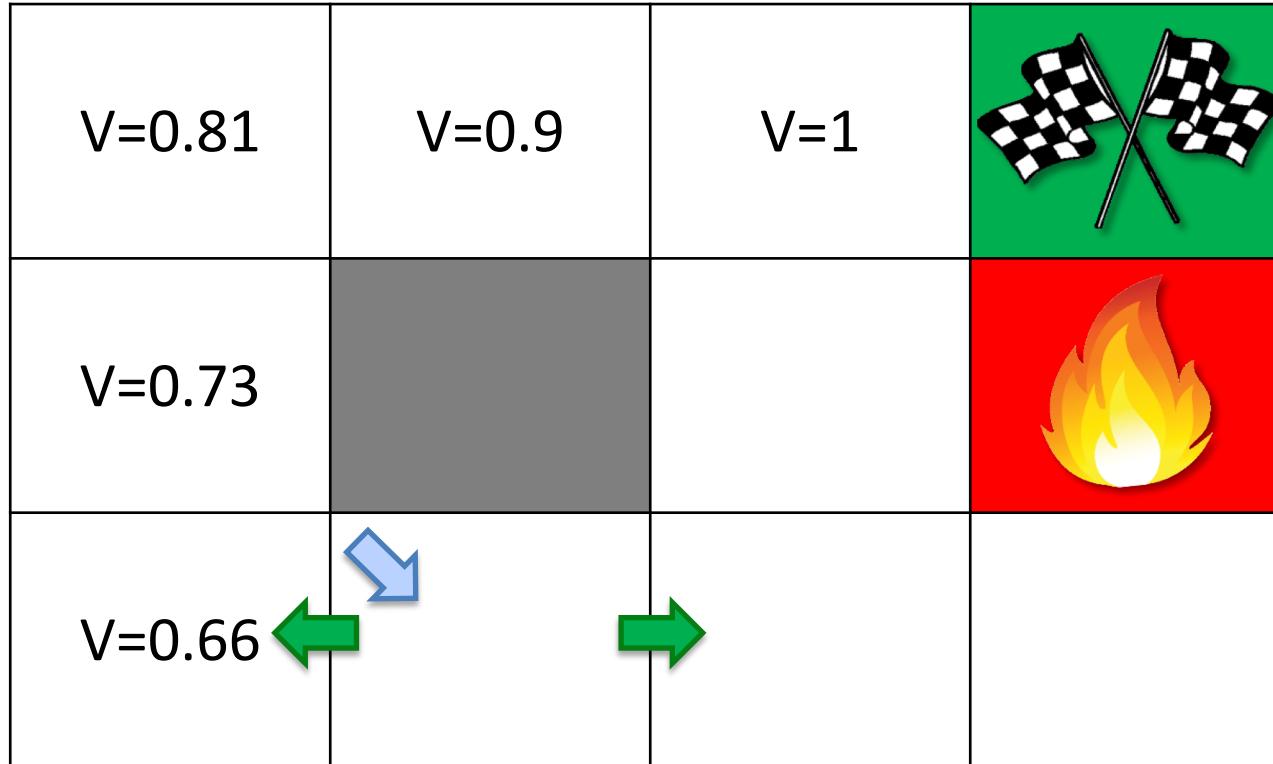
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73			
V=0.66			

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



$$\gamma = 0.9$$

A Equação de Bellman

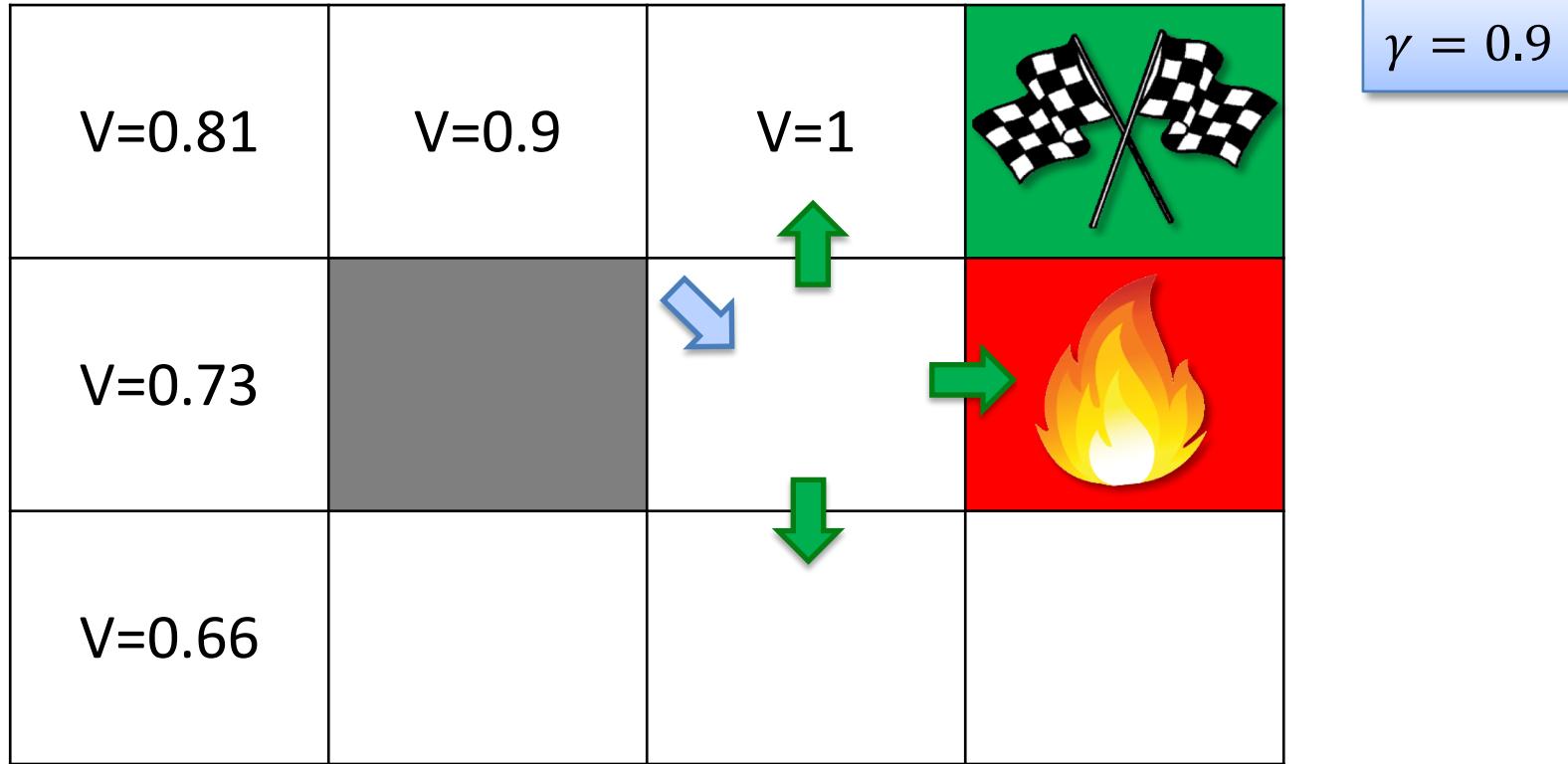
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73			
V=0.66			

$$\gamma = 0.9$$

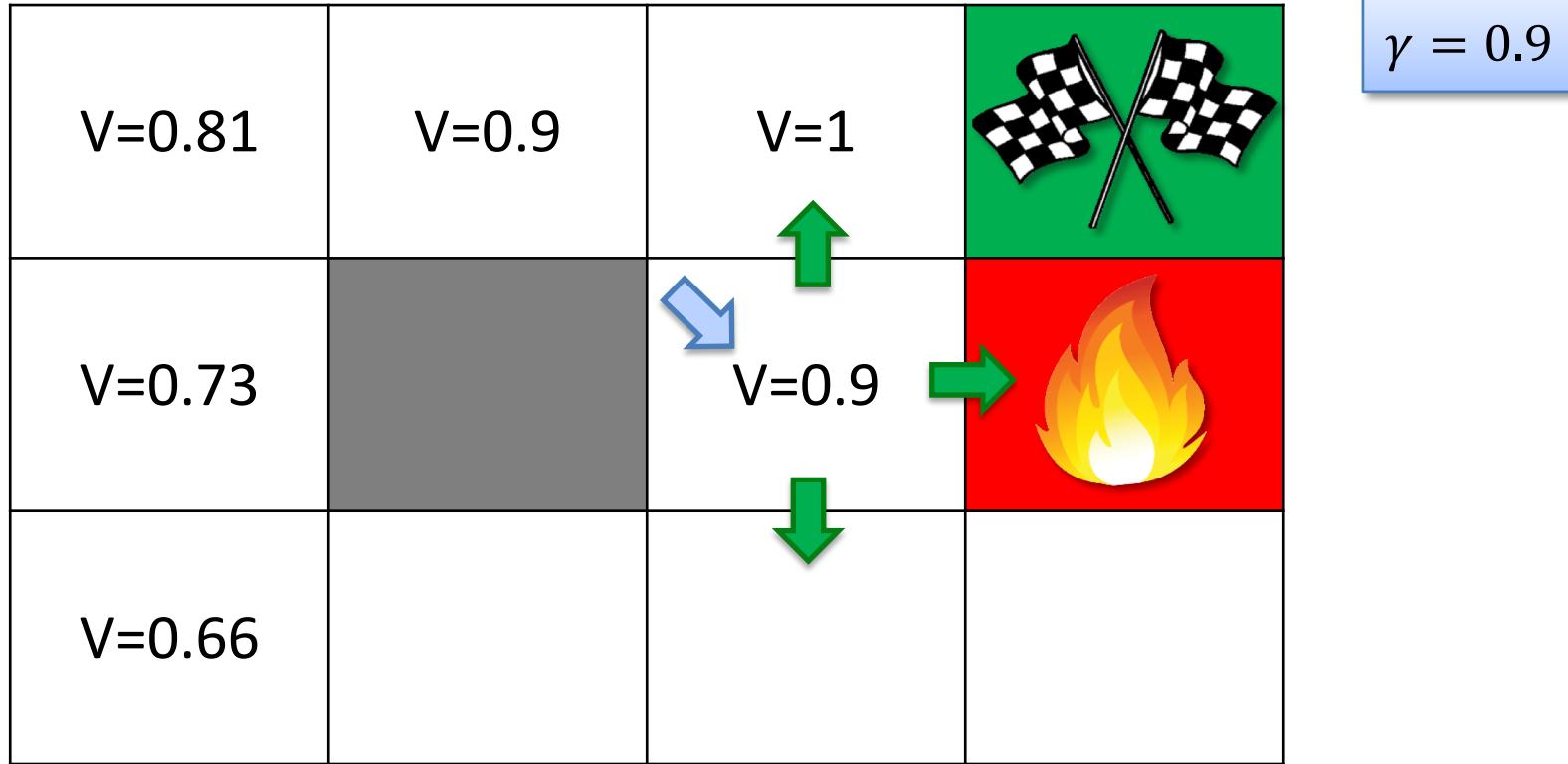
A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

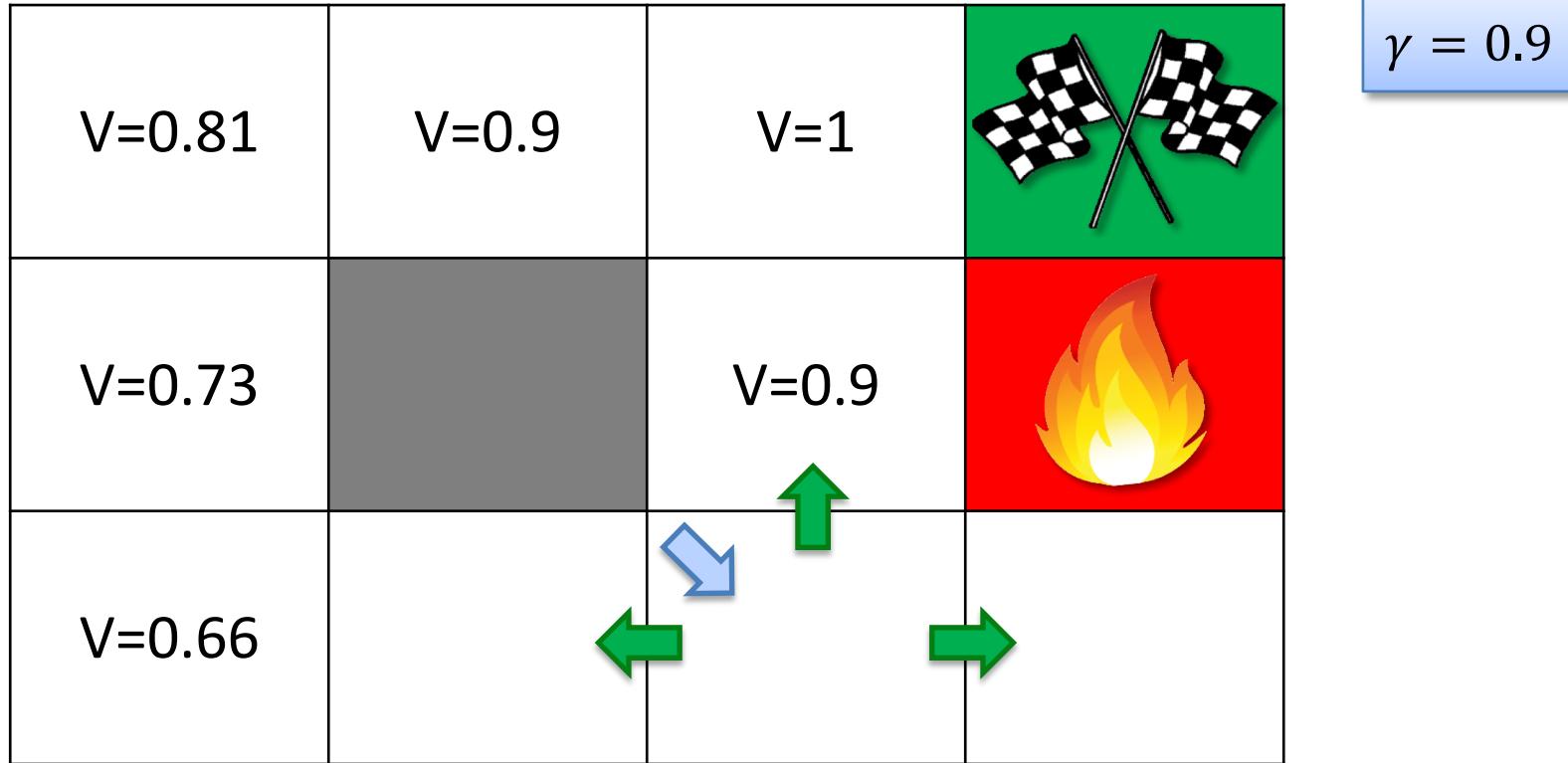
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73		V=0.9	
V=0.66			

$$\gamma = 0.9$$

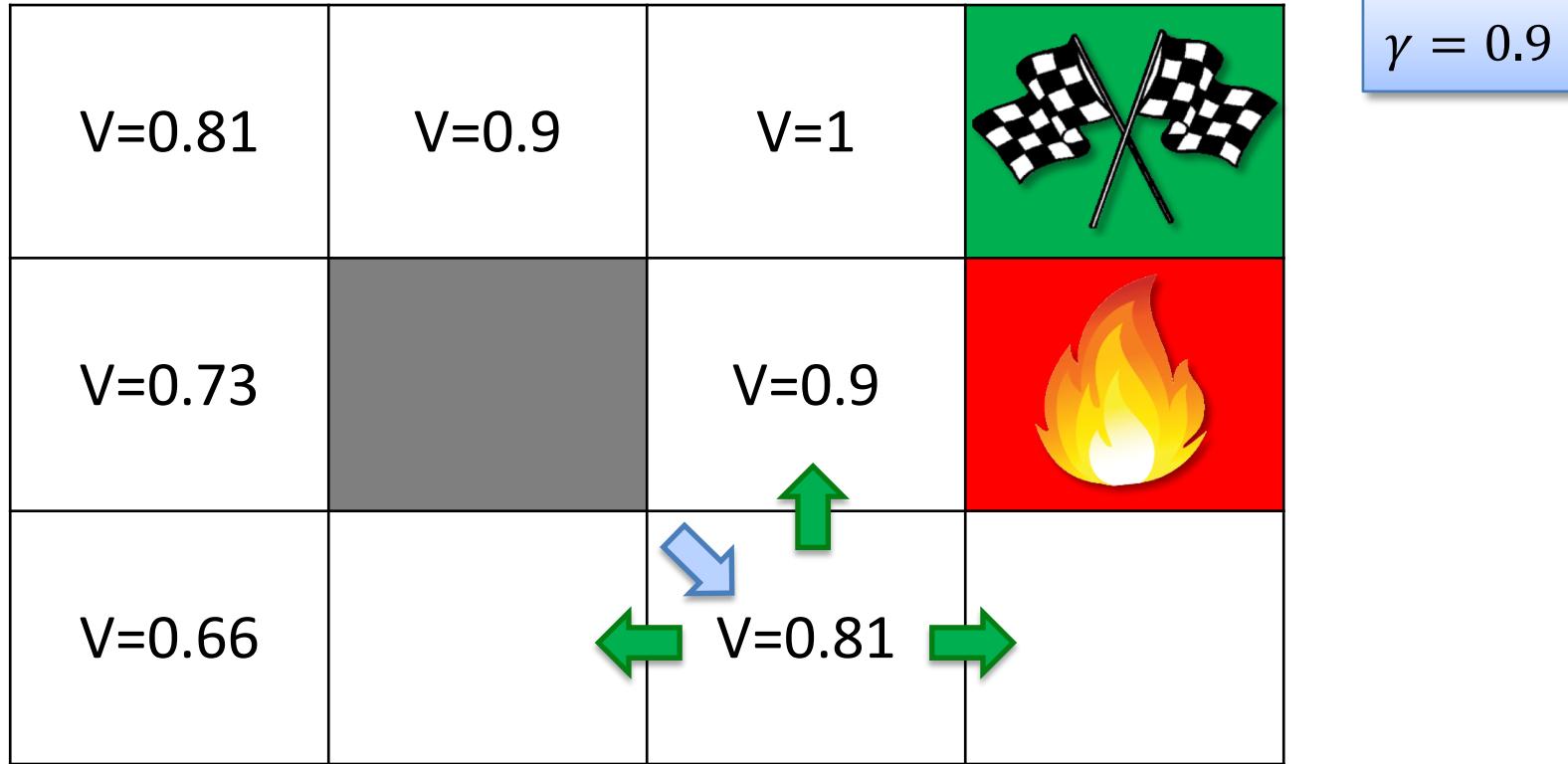
A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

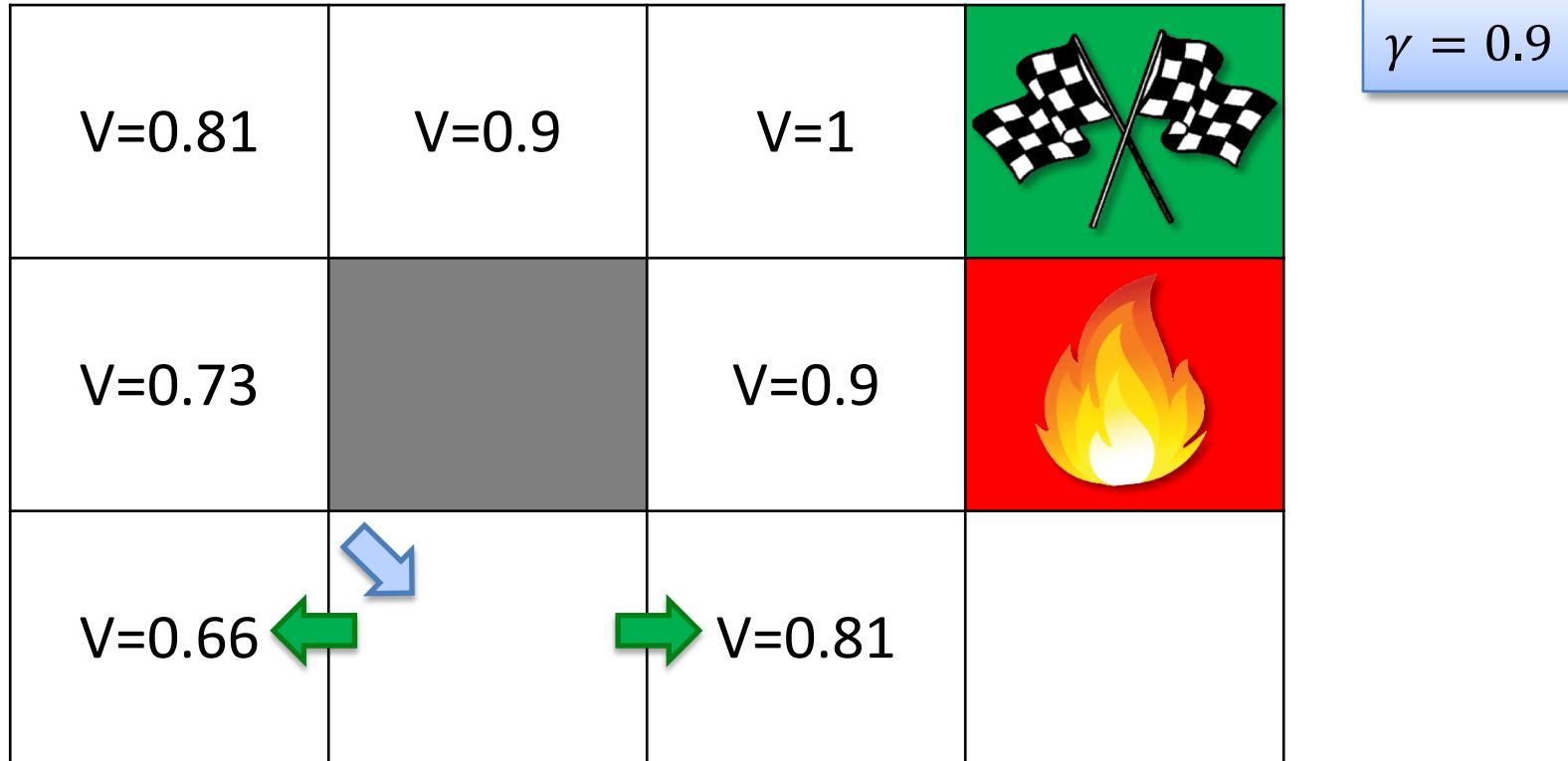
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73		V=0.9	
V=0.66		V=0.81	

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



A Equação de Bellman

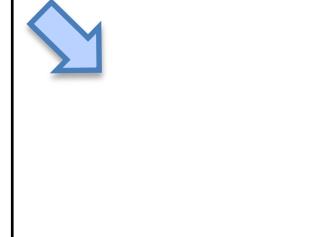
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1		
V=0.73		V=0.9		
V=0.66		V=0.73		V=0.81

$$\gamma = 0.9$$

A Equação de Bellman

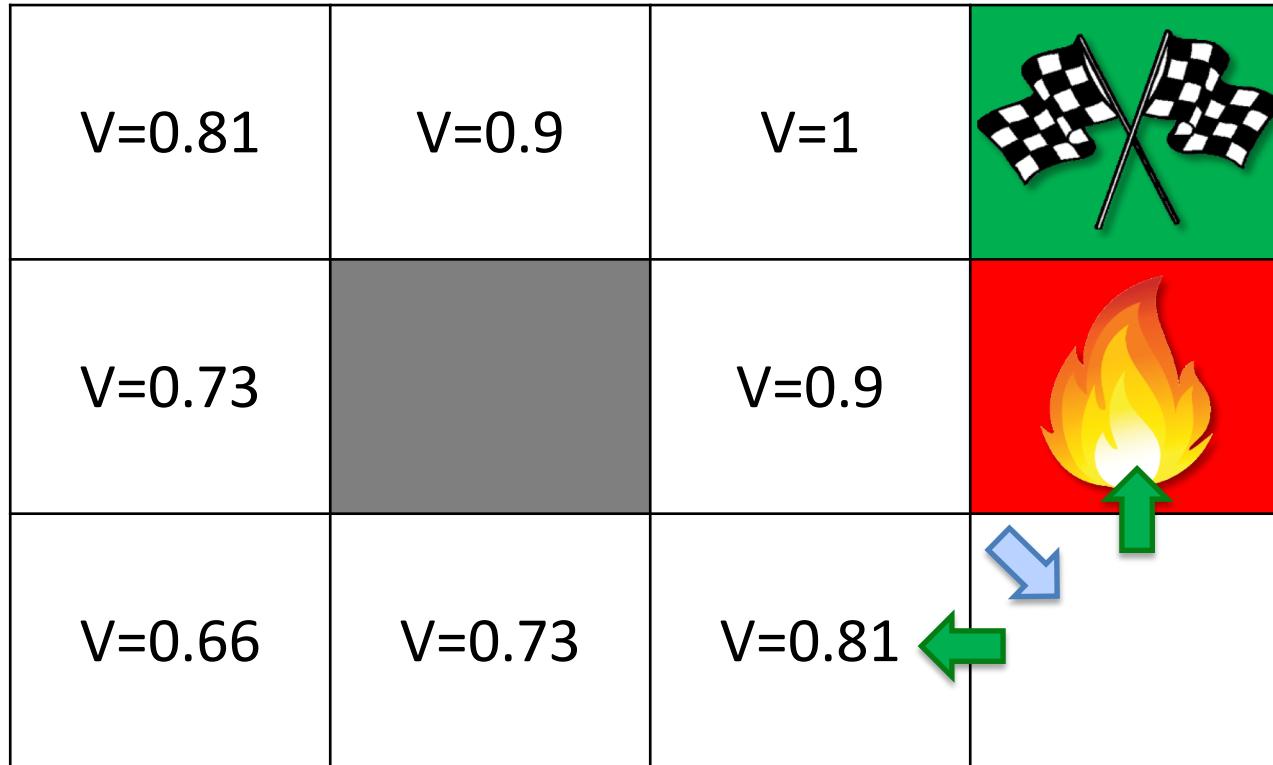
$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73		V=0.9	
V=0.66	V=0.73	V=0.81	

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$



$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73		V=0.9	
V=0.66	V=0.73	V=0.81	

$$\gamma = 0.9$$

A Equação de Bellman

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

V=0.81	V=0.9	V=1	
V=0.73		V=0.9	
V=0.66	V=0.73	V=0.81	V=0.73

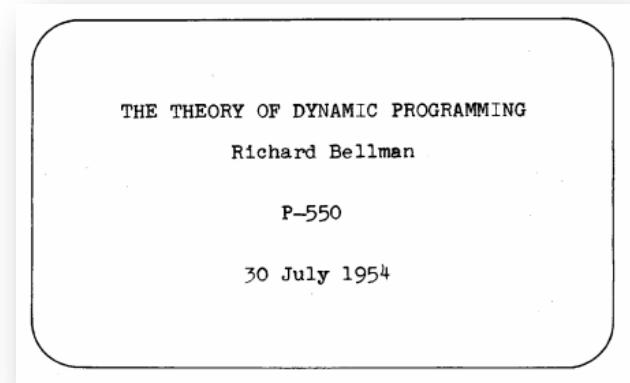
$$\gamma = 0.9$$

Leitura Adicional

Leitura Adicional:

The Theory of Dynamic Programming

Richard Bellman (1954)



Link:

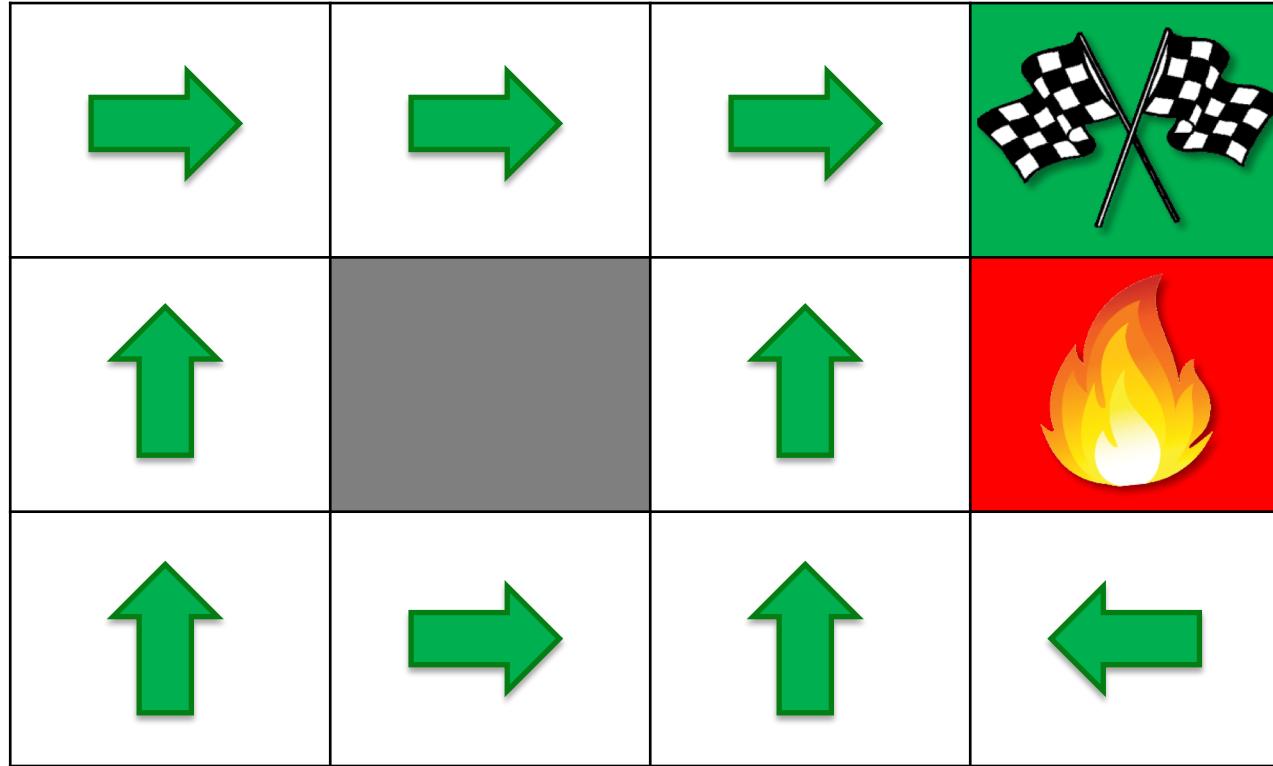
<https://www.rand.org/content/dam/rand/pubs/papers/2008/P550.pdf>

O “Plano”

A Equação de Bellman

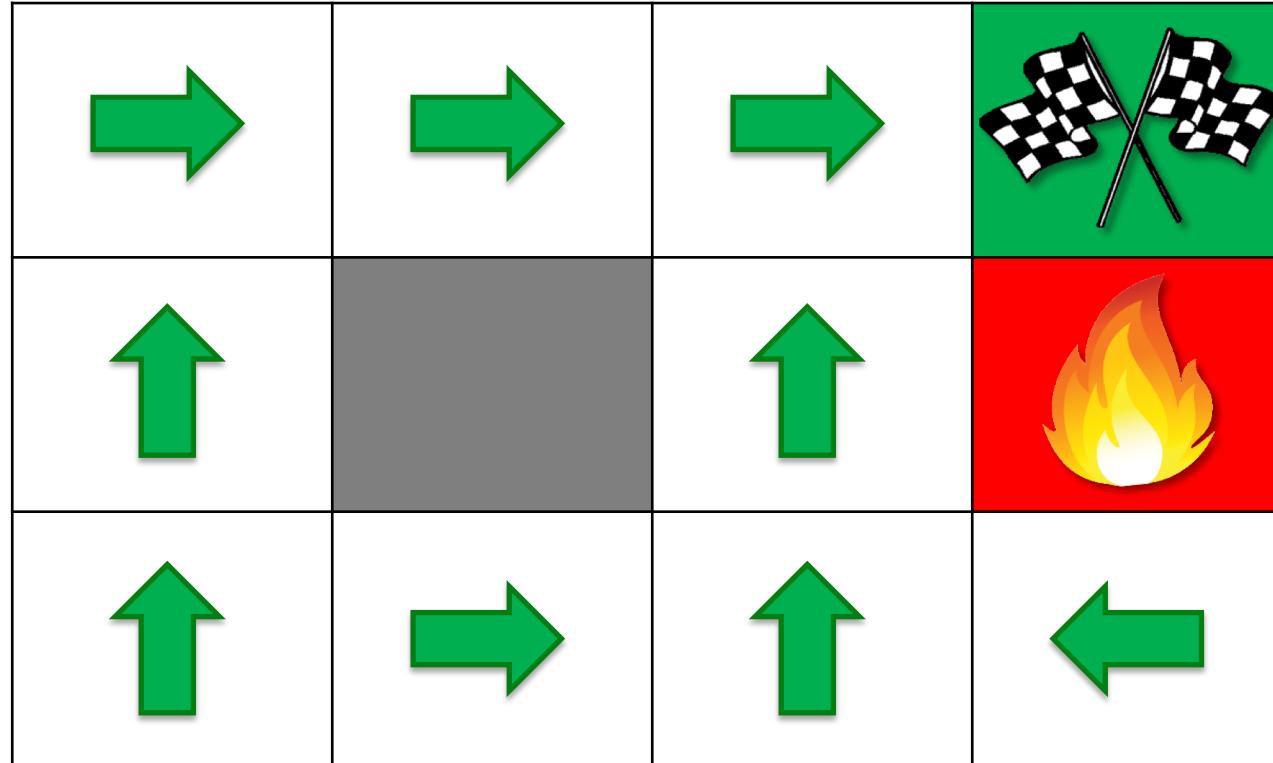
$V=0.81$	$V=0.9$	$V=1$	
$V=0.73$		$V=0.9$	
$V=0.66$	$V=0.73$	$V=0.81$	$V=0.73$

A Equação de Bellman

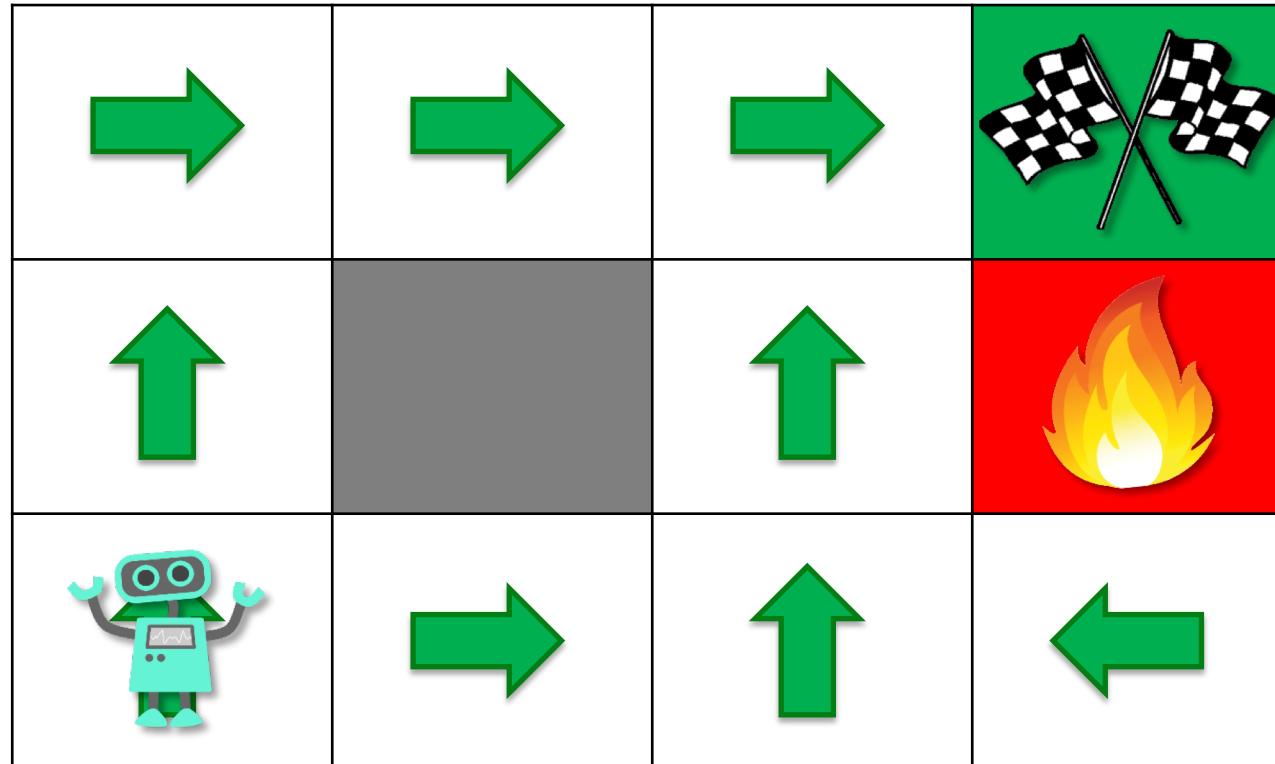


Markov Decision Process (MDP)

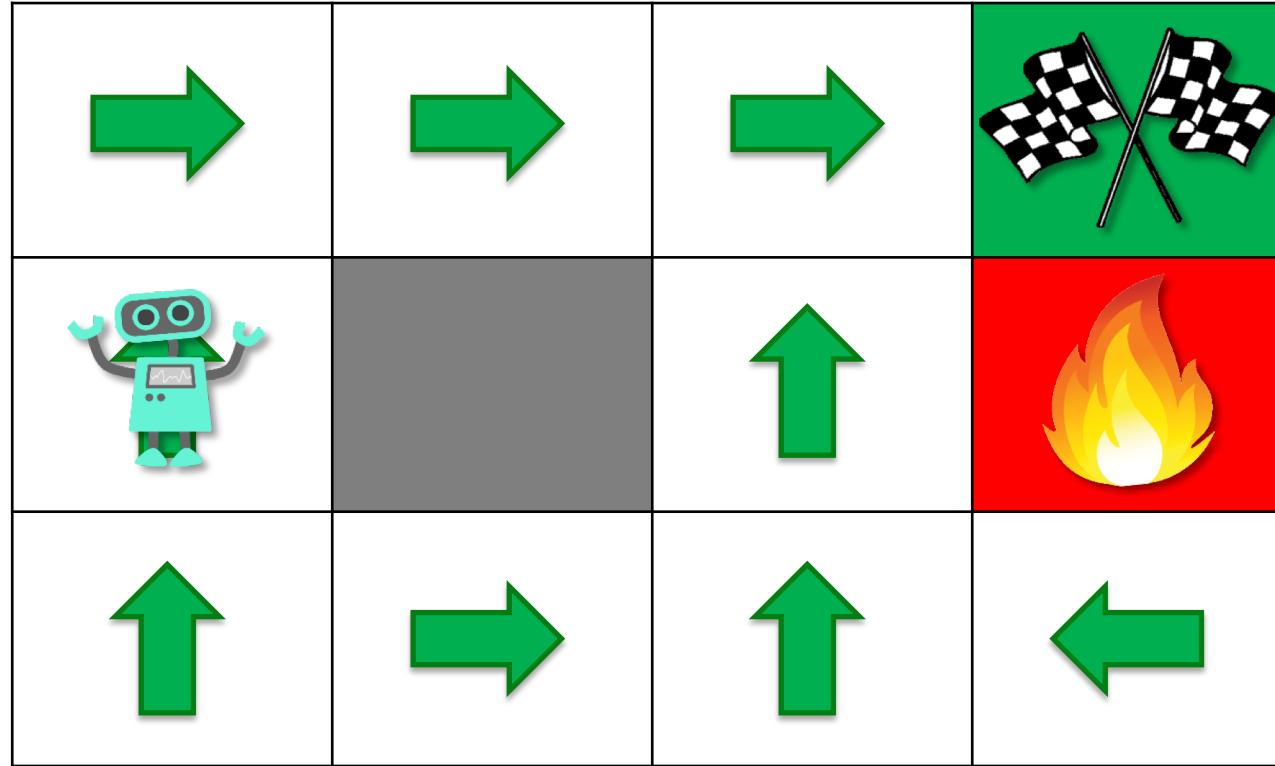
Markov Decision Process (MDP)



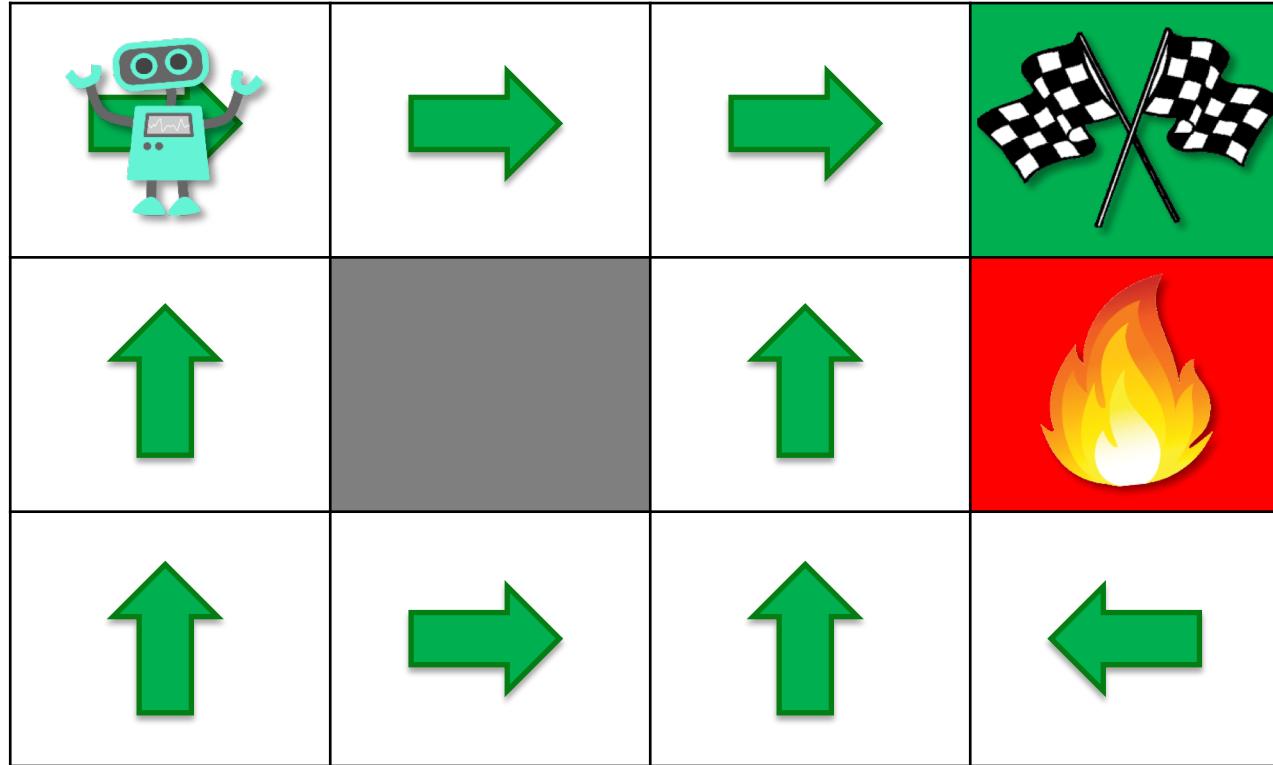
Markov Decision Process (MDP)



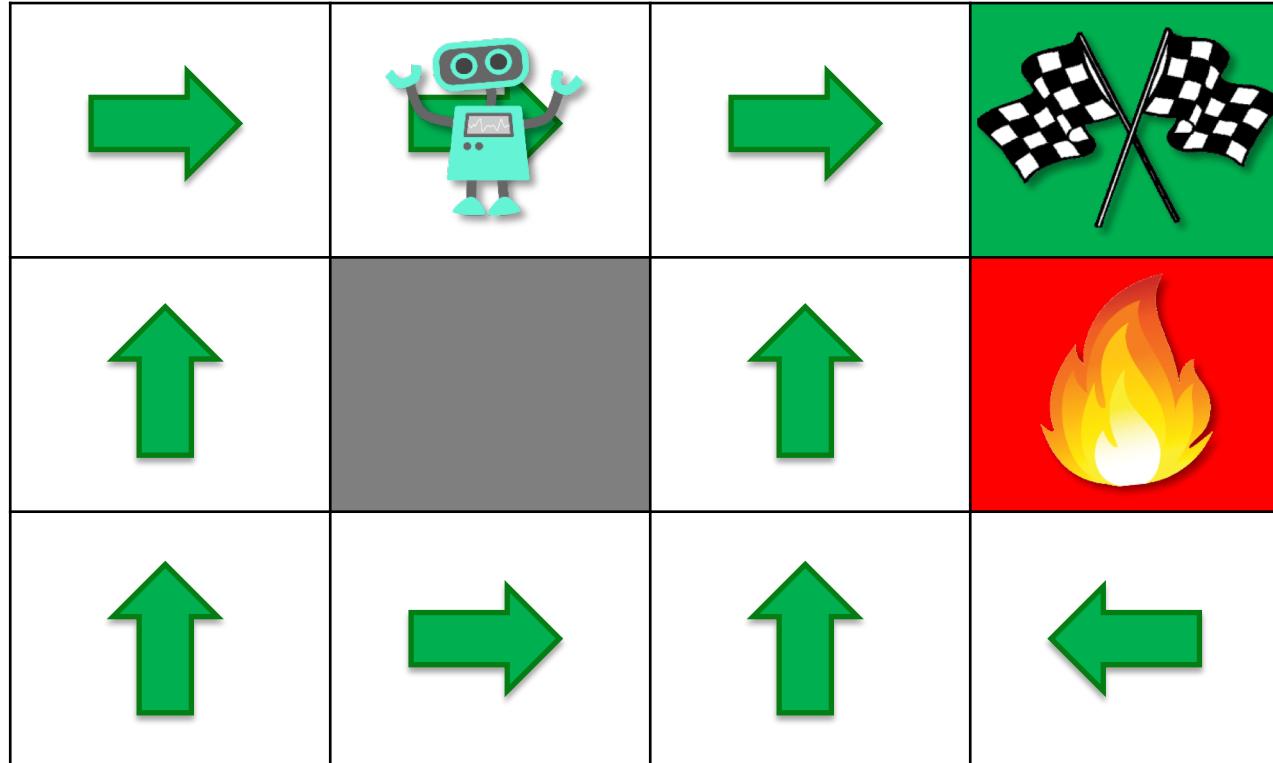
Markov Decision Process (MDP)



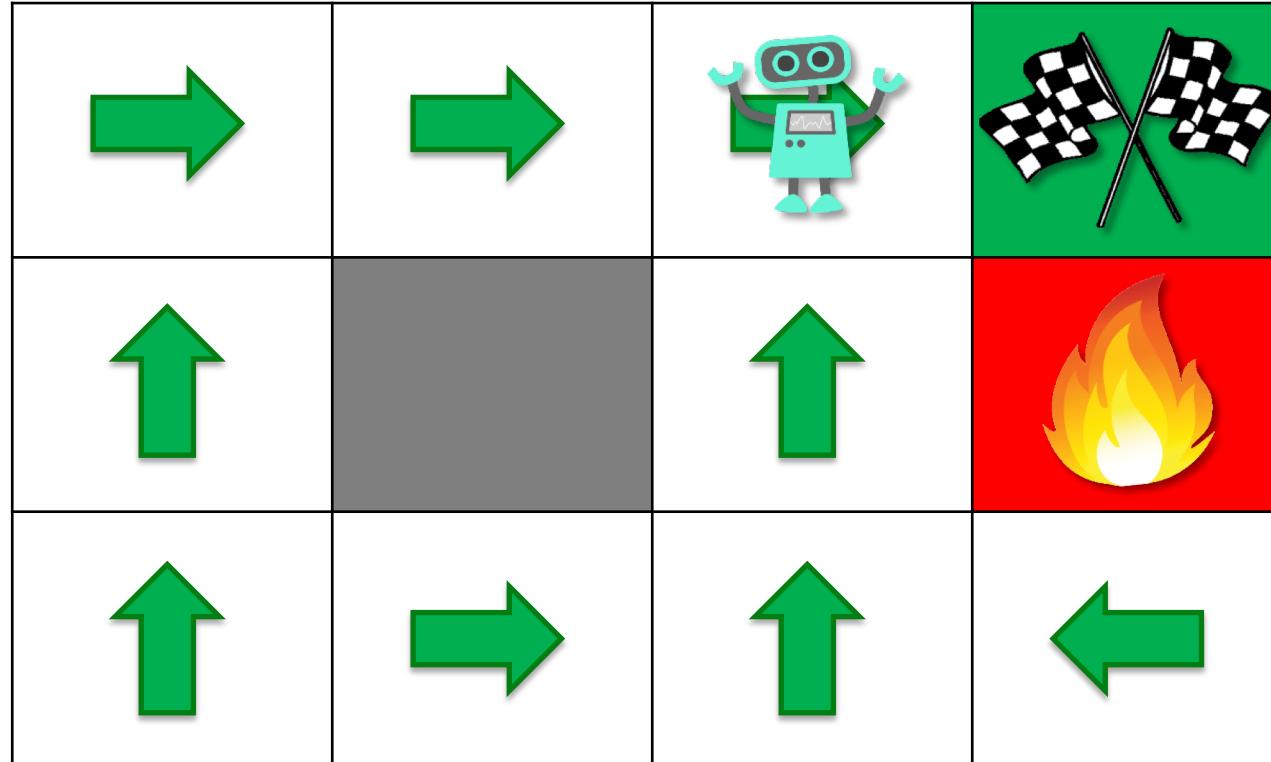
Markov Decision Process (MDP)



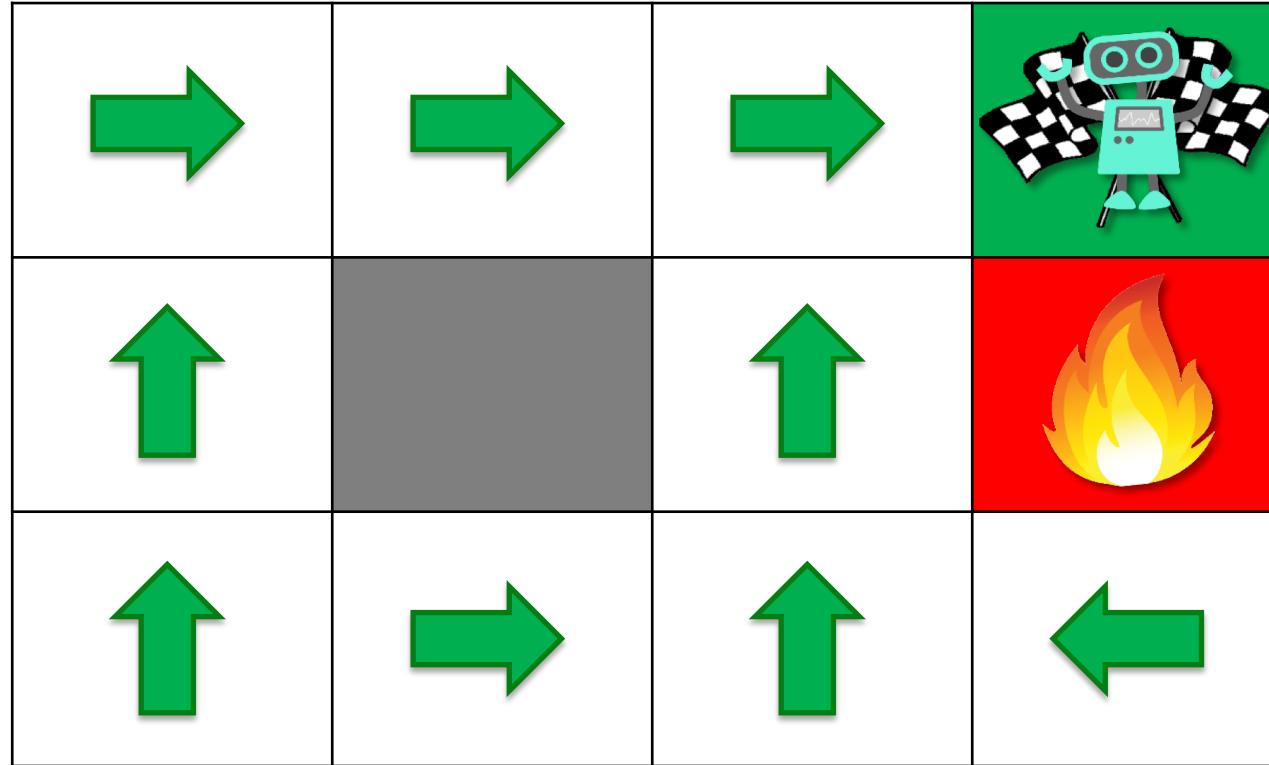
Markov Decision Process (MDP)



Markov Decision Process (MDP)

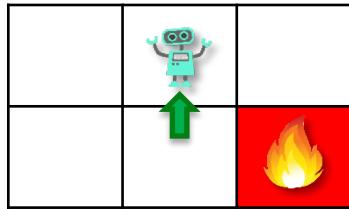
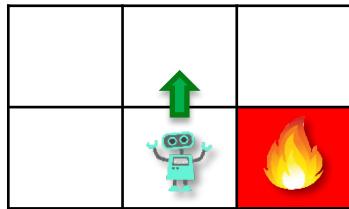


Markov Decision Process (MDP)

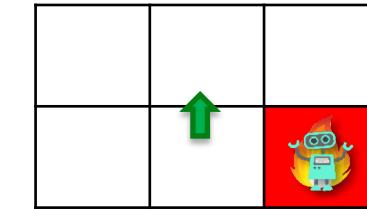
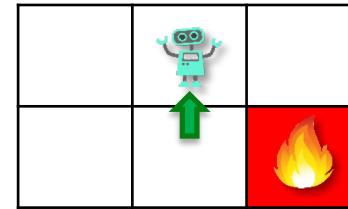
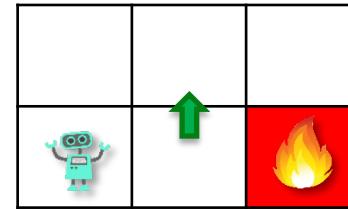
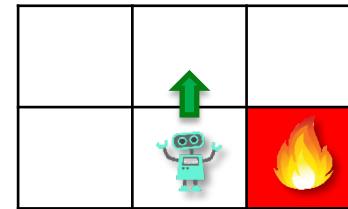


Markov Decision Process (MDP)

Busca Determinística



Busca Não-Determinística



Markov Decision Process (MDP)

Markov Process

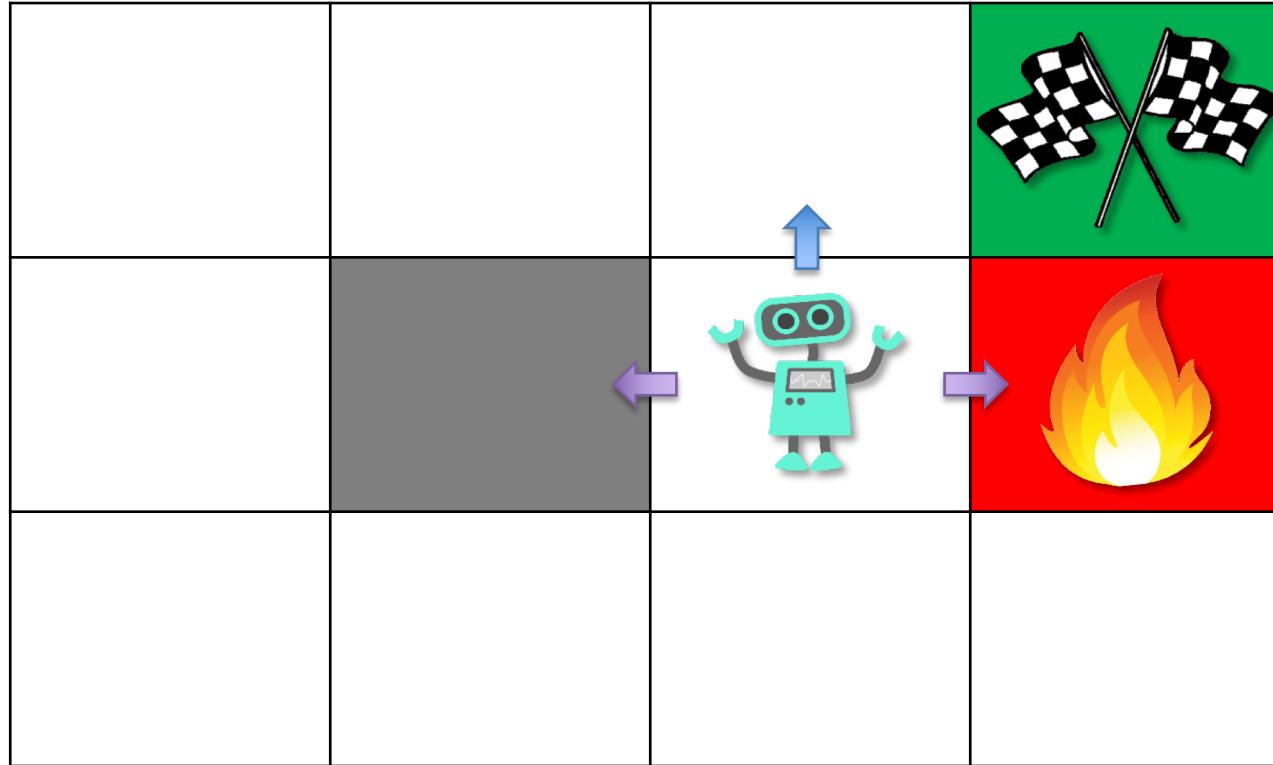
Markov Decision Process (MDP)

Markov Decision Process (MDP)

Um processo estocástico possui a propriedade Markov se a distribuição de probabilidade condicional de estados futuros do processo (condicional em estados passados e presentes) depende apenas do estado presente, não da sequência de eventos que o precedeu. Um processo com essa propriedade é chamado de processo de Markov.

Wikipedia (traduzido para português)

Markov Decision Process (MDP)

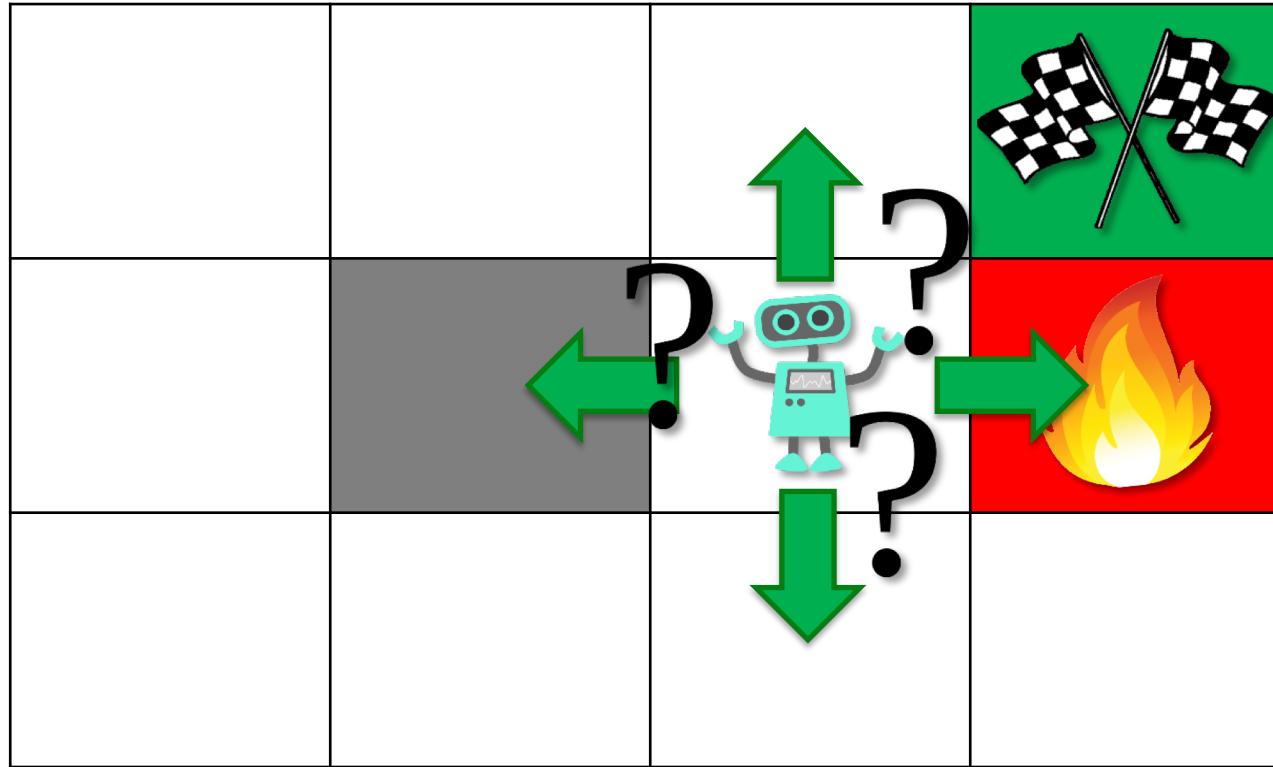


Markov Decision Process (MDP)

Markov Decision Processes (MDPs – Processos de Decisão de Makov) fornecem uma estrutura matemática para modelar a tomada de decisões em situações onde os resultados são parcialmente aleatórios e parcialmente sob o controle de um tomador de decisão.

Wikipedia (traduzido para o português)

Markov Decision Process (MDP)



Markov Decision Process (MDP)

$$V(s) = \max_a(R(s, a) + \gamma V(s'))$$

Markov Decision Process (MDP)

$$V(s) = \max_a (R(s, a) + \gamma \overbrace{V(s')}^{\text{Red bracket}})$$

Markov Decision Process (MDP)

$$V(s) = \max_a (R(s, a) + \gamma \overbrace{V(s')}^{\substack{s'_1 & s'_2 & s'_3}})$$

Markov Decision Process (MDP)

$$V(s) = \max_a (R(s, a) + \gamma \overbrace{V(s')}^{\text{Red bracket}})$$

$V(s'_1) \quad V(s'_2) \quad V(s'_3)$

Markov Decision Process (MDP)

$$V(s) = \max_a (R(s, a) + \gamma \overbrace{V(s')}^{\text{Red bracket}})$$

$0.8 * V(s'_1) + 0.1 * V(s'_2) + 0.1 * V(s'_3)$

Markov Decision Process (MDP)

$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

Leitura Adicional

Leitura Adicional:

A Survey of Applications of Markov Decision Processes

D. J. White (1993)

Link:

<http://www.cs.uml.edu/ecg/uploads/AIfall14/MDPApplications3.pdf>

TABLE 1. Application areas

1	Population harvesting	(5)
2	Agriculture	(5)
3	Water resources	(15)
4	Inspection, maintenance and repair	(18)
5	Purchasing, inventory and production	(14)
6	Finance and investment	(9)
7	Queues	(6)
8	Sales promotion	(4)
9	Search	(3)
10	Motor insurance claims	(2)
11	Overbooking	(5)
12	Epidemics	(2)
13	Credit	(2)
14	Sports	(2)
15	Patient admissions	(1)
16	Location	(1)
17	Design of experiments	(1)
18	General	(5)

Política x Plano

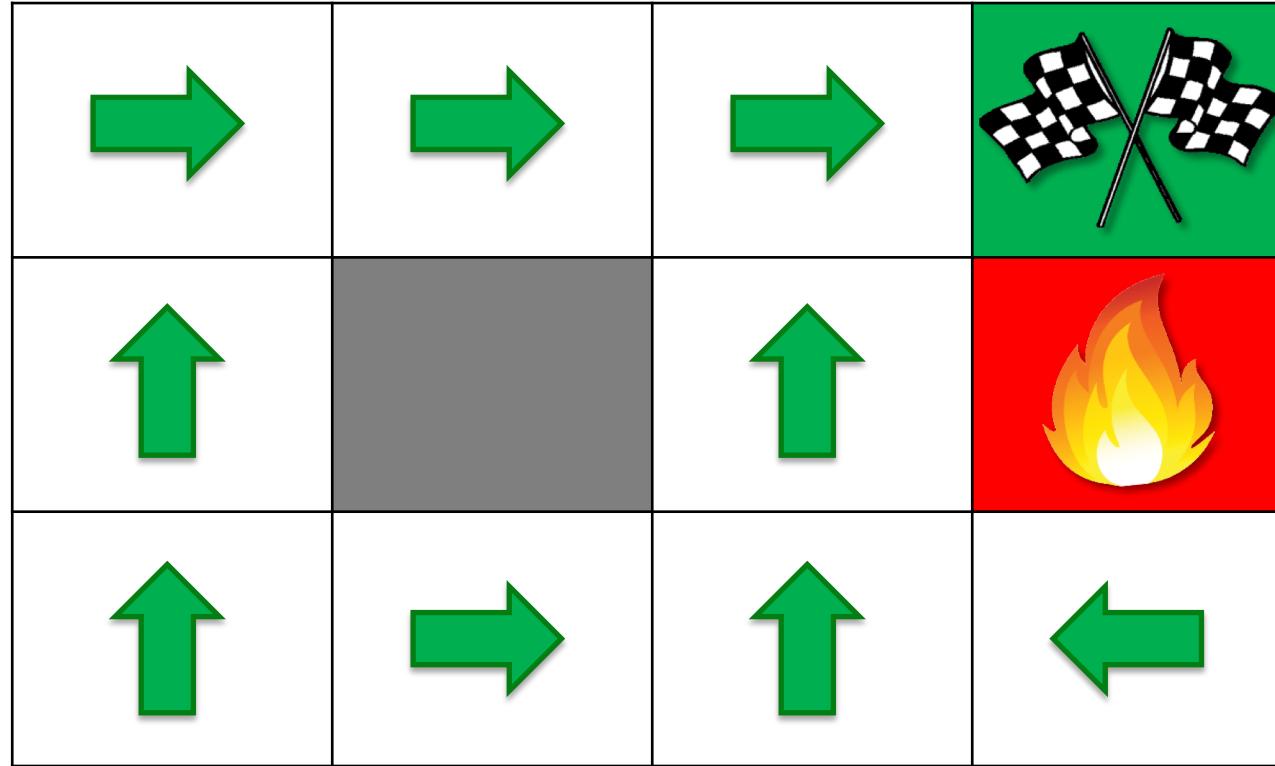
Política x Plano

$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

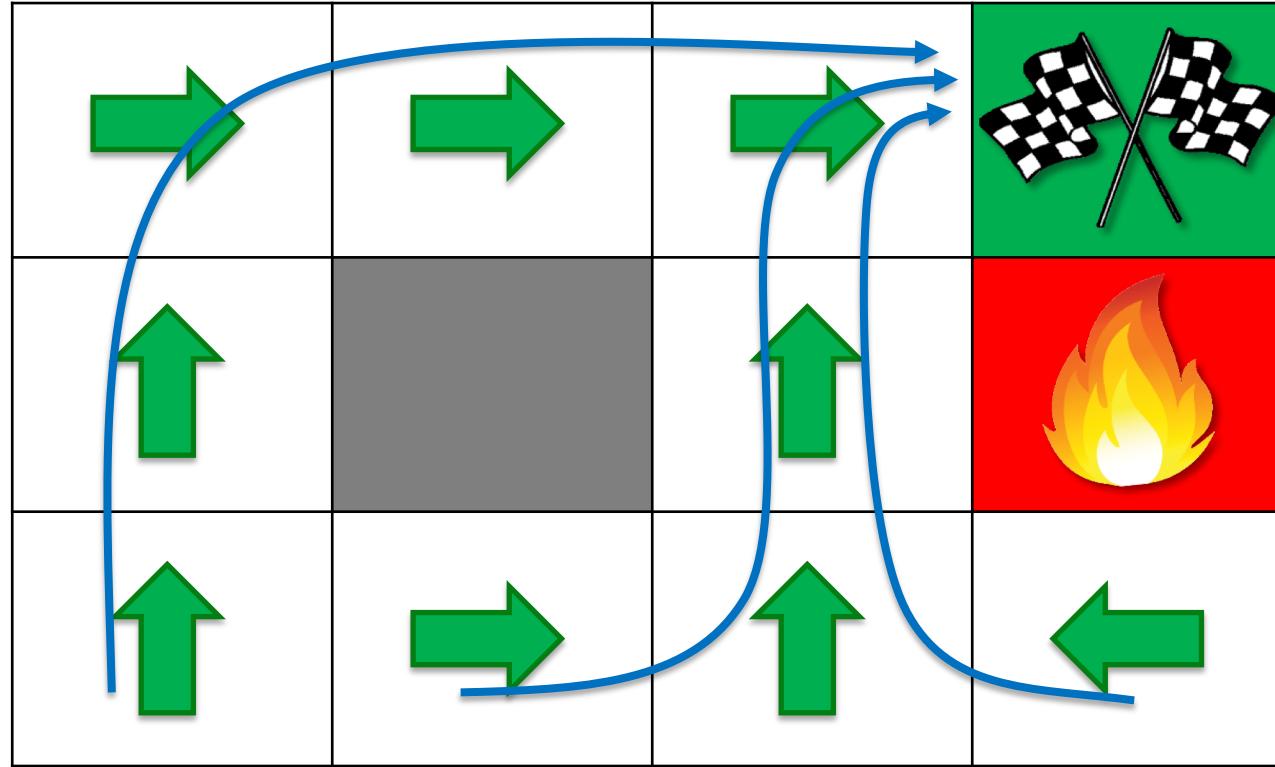
Política x Plano

$V=0.81$	$V=0.9$	$V=1$	
$V=0.73$		$V=0.9$	
$V=0.66$	$V=0.73$	$V=0.81$	$V=0.73$

Política x Plano



Política x Plano



Política x Plano

$V=0.71$	$V=0.74$	$V=0.86$	
$V=0.63$		$V=0.39$	
$V=0.55$	$V=0.46$	$V=0.36$	$V=0.22$

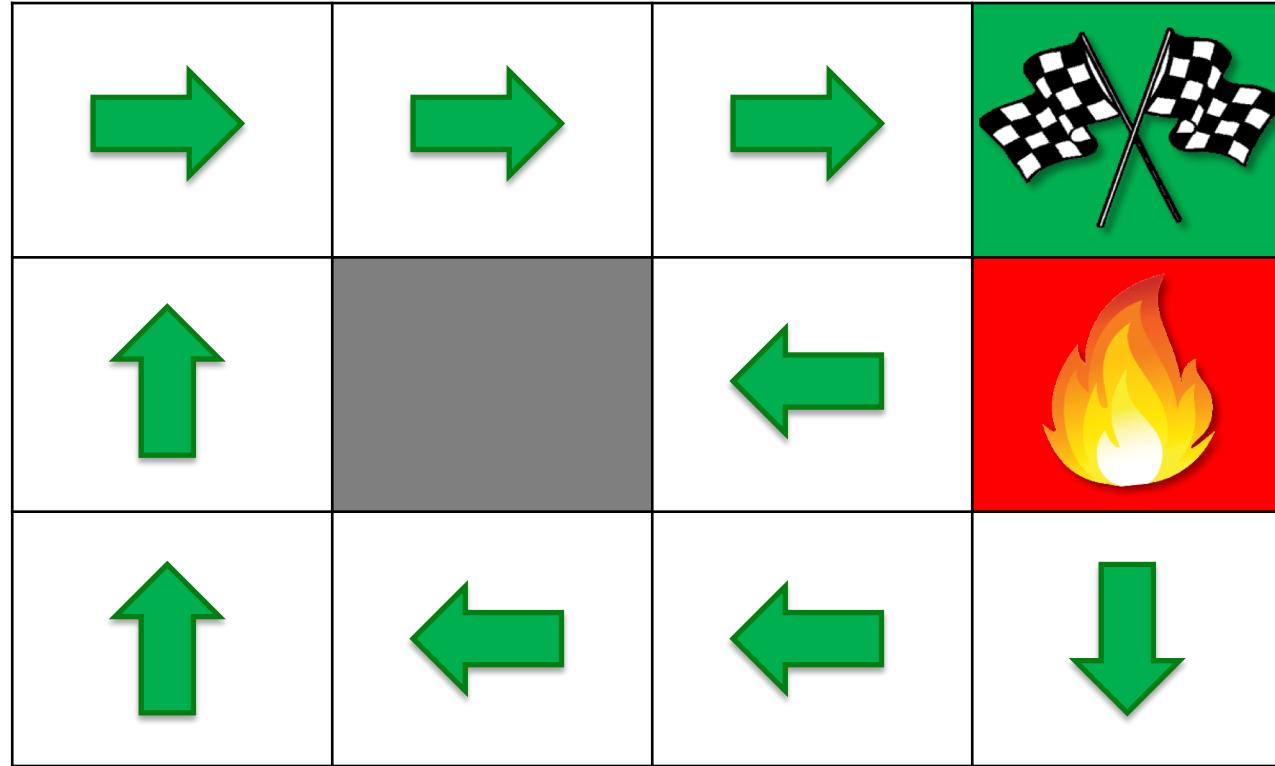
Política x Plano

$V=0.81$	$V=0.9$	$V=1$	
$V=0.73$		$V=0.9$	
$V=0.66$	$V=0.73$	$V=0.81$	$V=0.73$

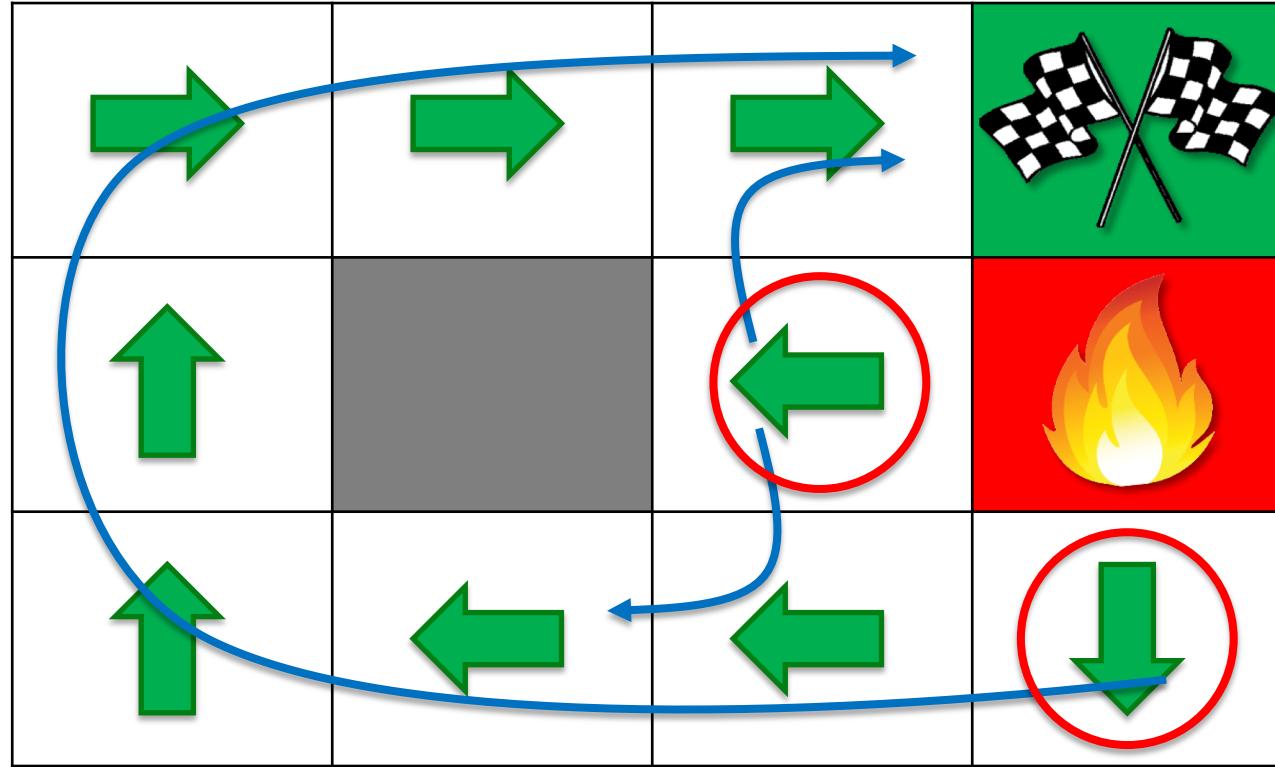
Política x Plano

$V=0.71$	$V=0.74$	$V=0.86$	
$V=0.63$		$V=0.39$	
$V=0.55$	$V=0.46$	$V=0.36$	$V=0.22$

Política x Plano



Política x Plano



Adição de Penalidades – “Living Penalty”

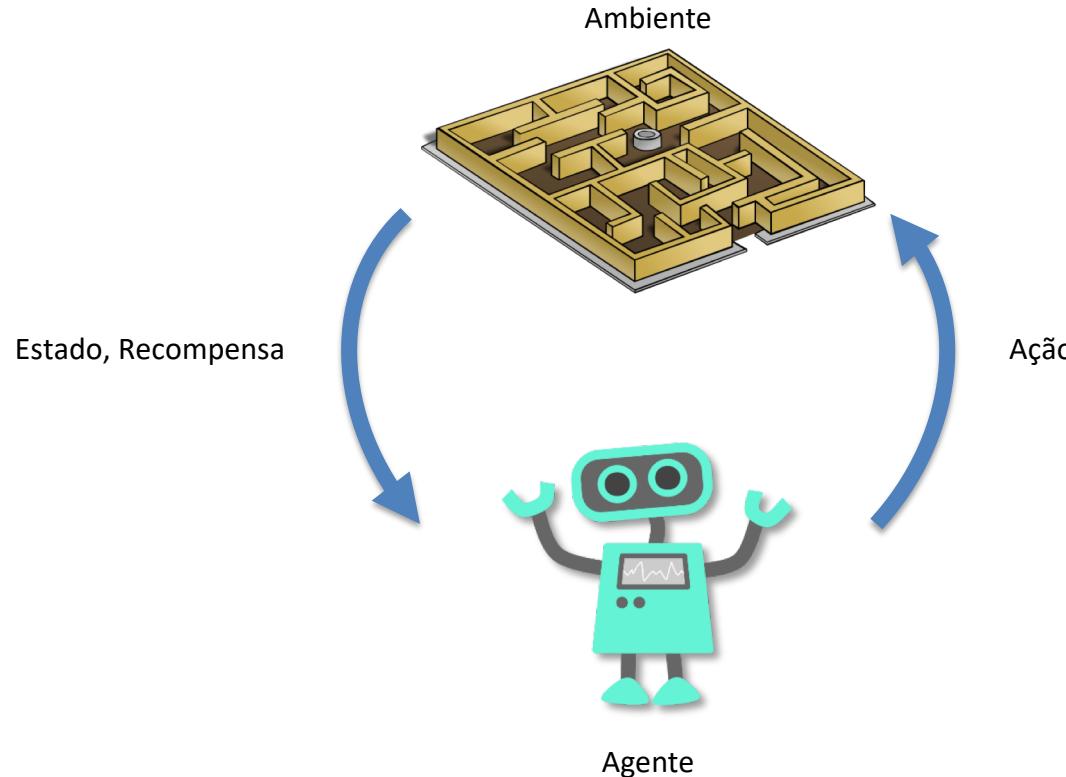
“Living Penalty”

$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

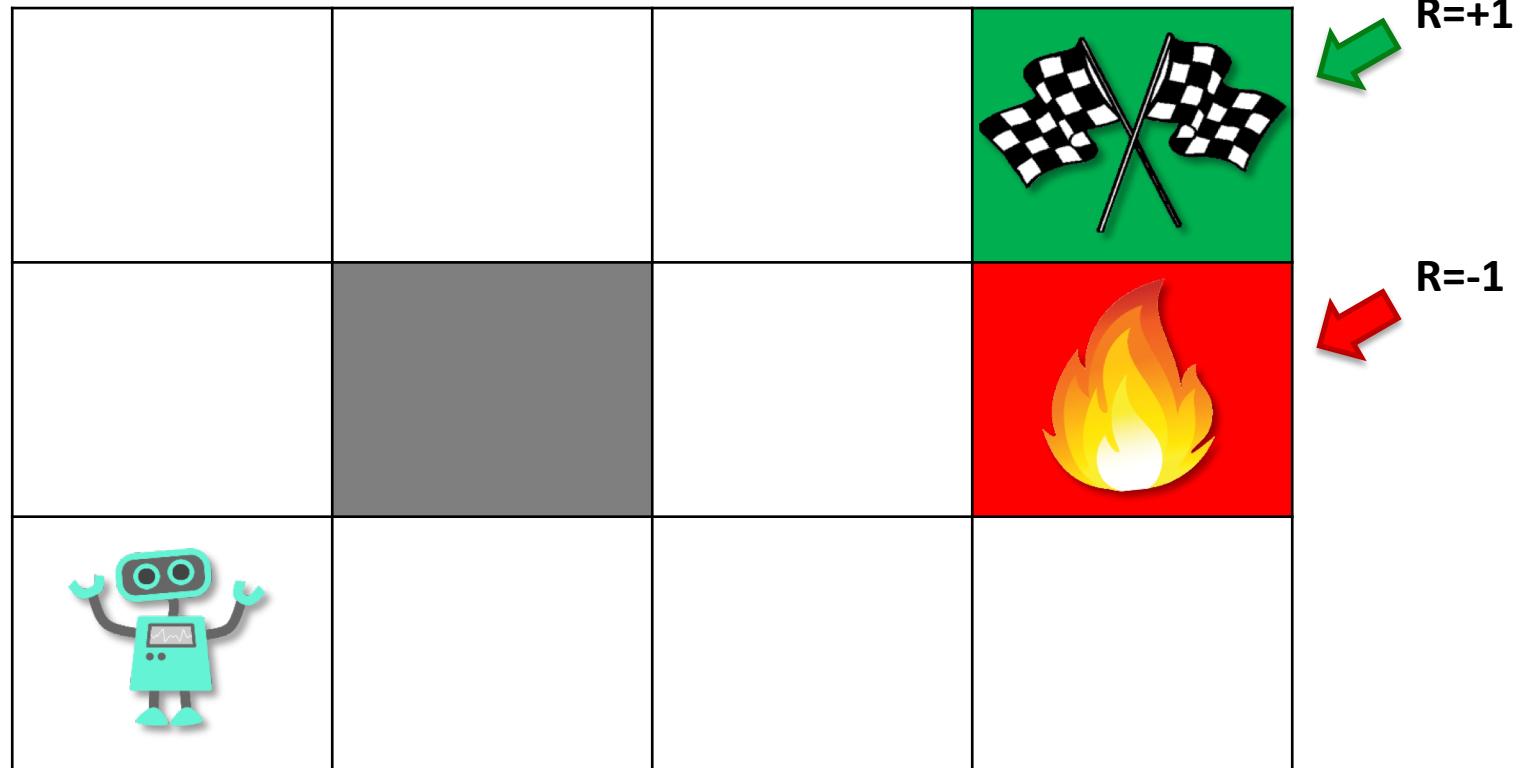
“Living Penalty”

$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

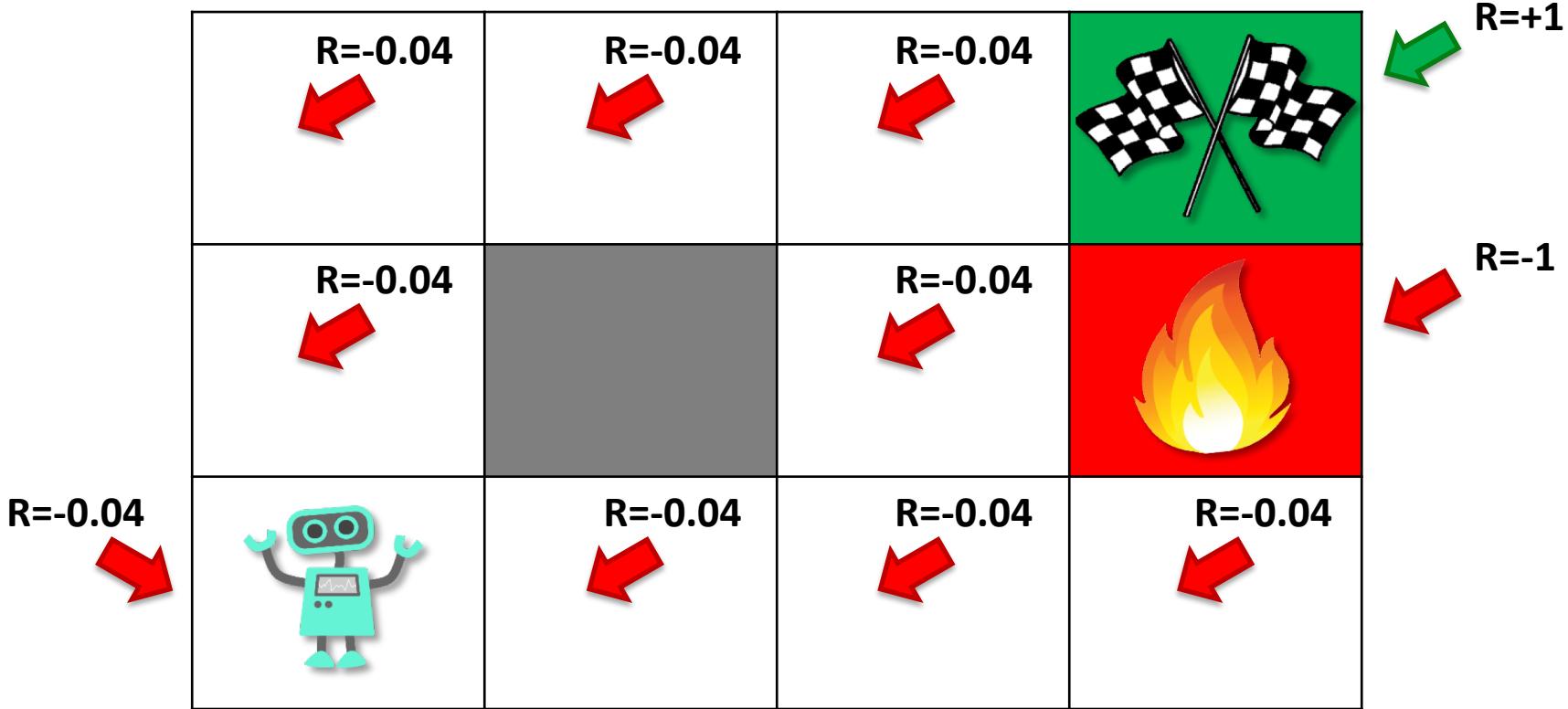

“Living Penalty”



“Living Penalty”

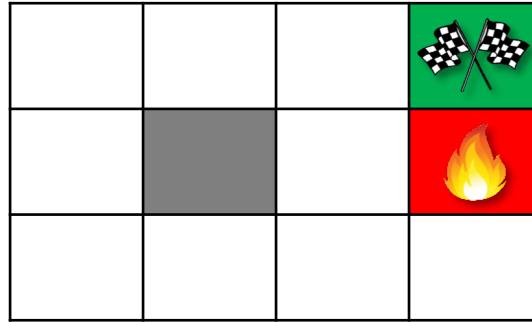


“Living Penalty”

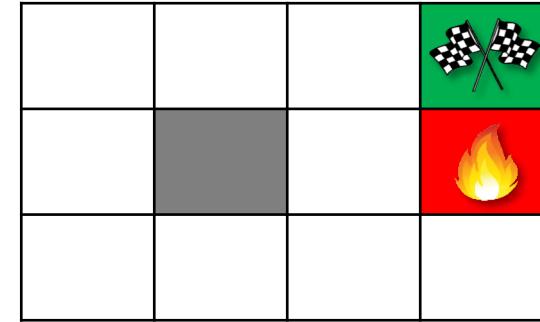


“Living Penalty”

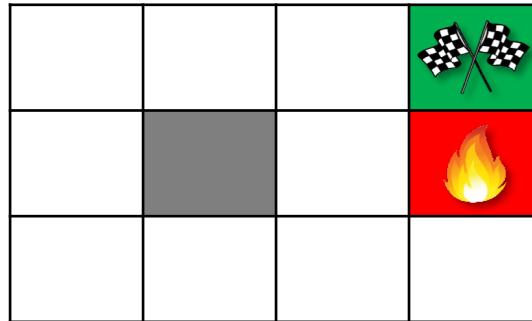
$R(s)=0$



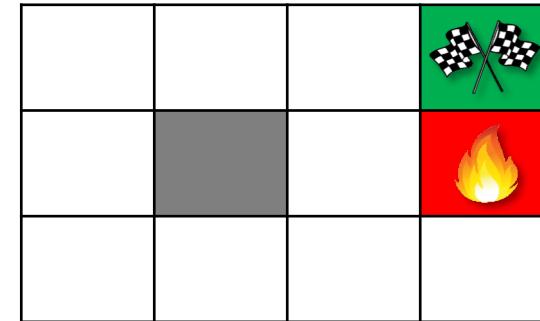
$R(s)=-0.04$



$R(s)=-0.5$

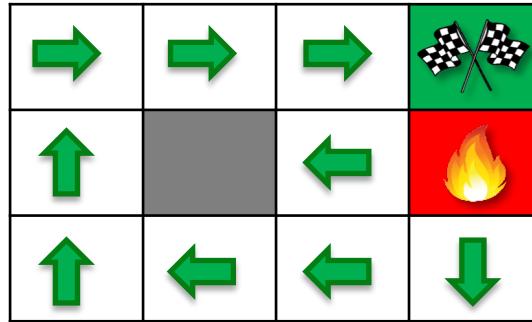


$R(s)=-2.0$

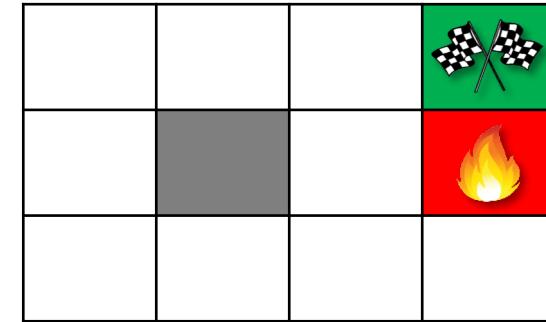


“Living Penalty”

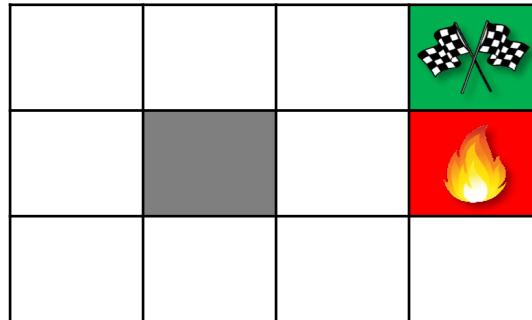
$R(s)=0$



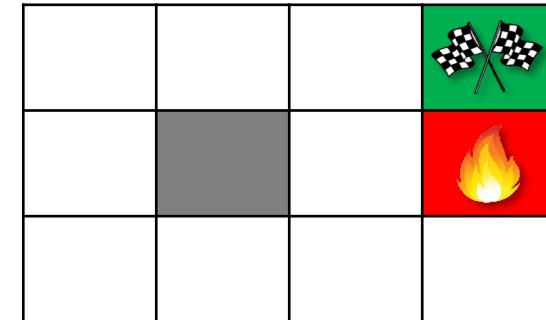
$R(s)=-0.04$



$R(s)=-0.5$

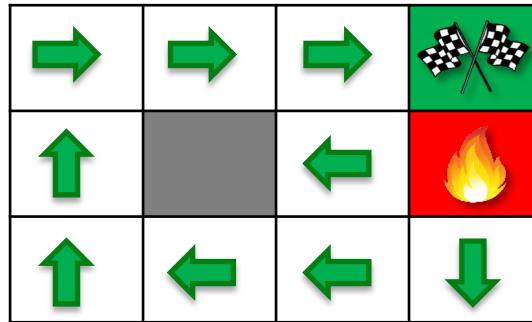


$R(s)=-2.0$

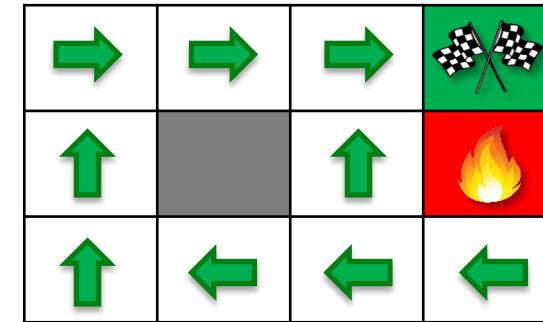


“Living Penalty”

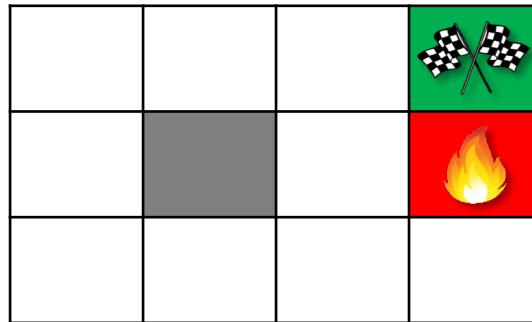
$R(s)=0$



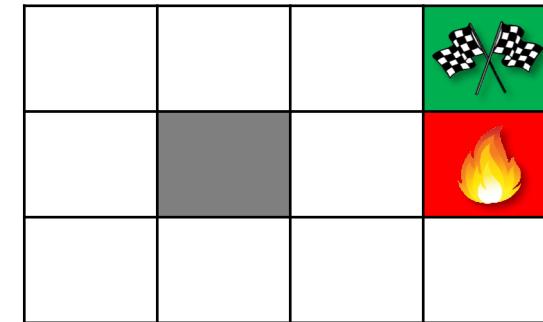
$R(s)=-0.04$



$R(s)=-0.5$

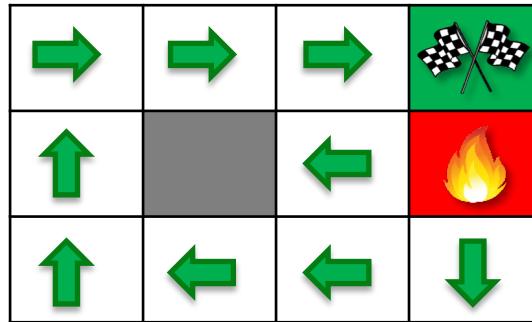


$R(s)=-2.0$

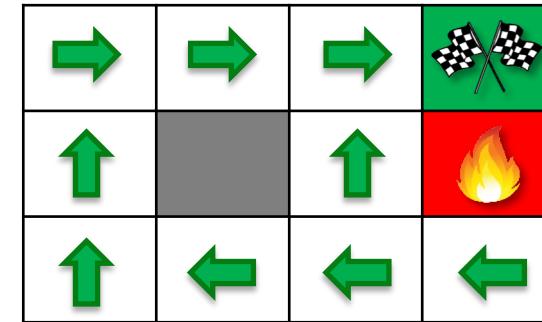


“Living Penalty”

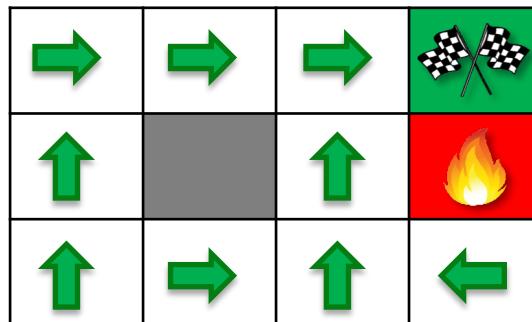
$R(s)=0$



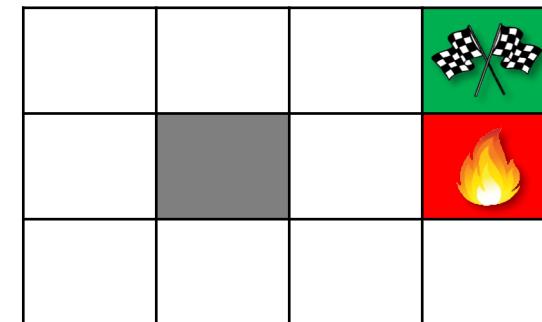
$R(s)=-0.04$



$R(s)=-0.5$

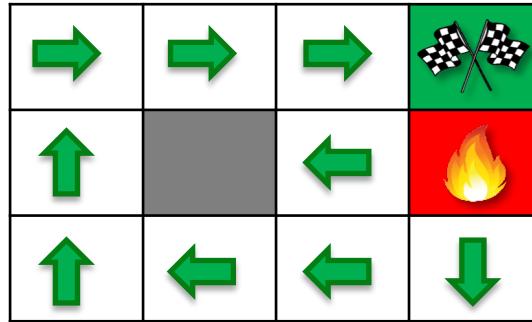


$R(s)=-2.0$

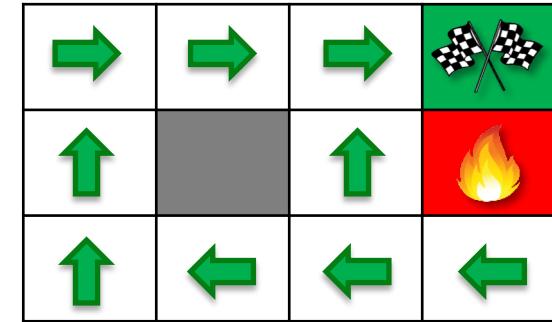


“Living Penalty”

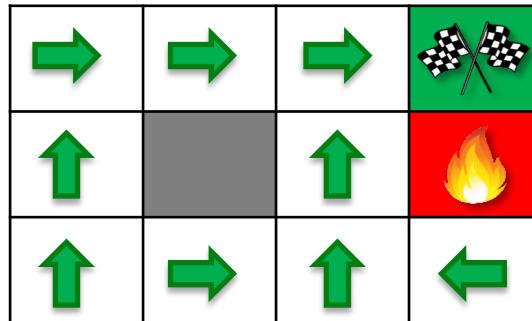
$R(s)=0$



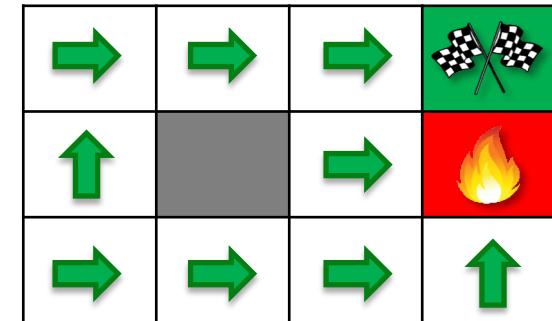
$R(s)=-0.04$



$R(s)=-0.5$



$R(s)=-2.0$



Q-Learning Intuição

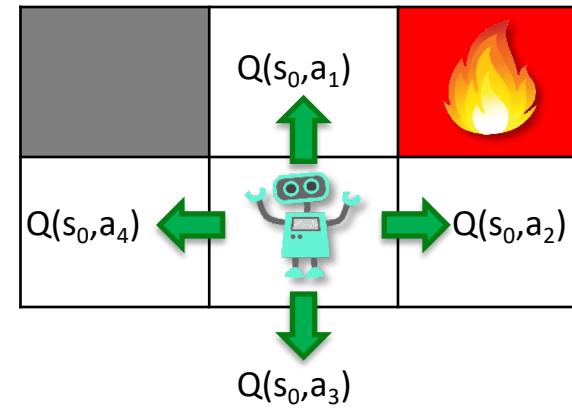
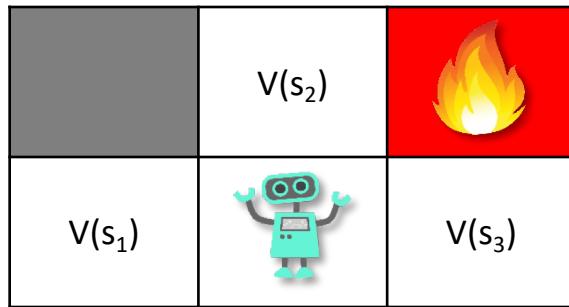
Q-Learning Intuição

$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

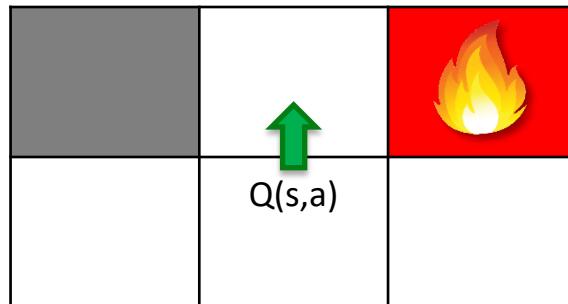
Q-Learning Intuição

Onde está o Q?

Q-Learning Intuição



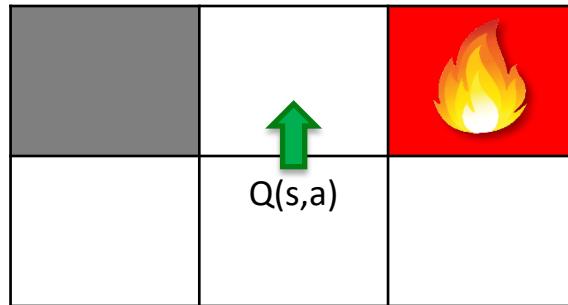
Q-Learning Intuição



Q-Learning Intuição



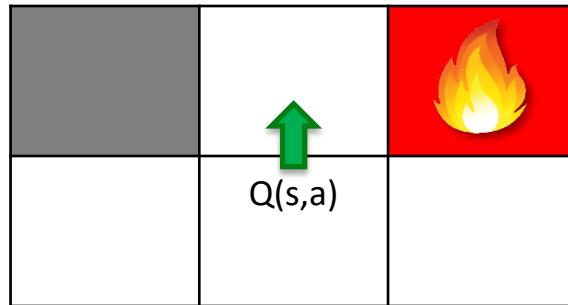
$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$



Q-Learning Intuição



$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

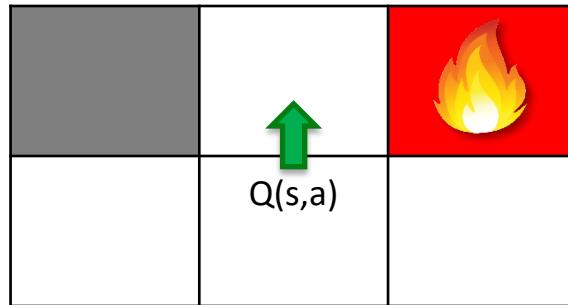


$$Q(s, a) =$$

Q-Learning Intuição



$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

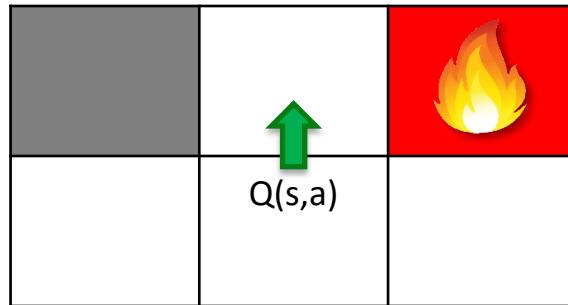


$$Q(s, a) = R(s, a) +$$

Q-Learning Intuição

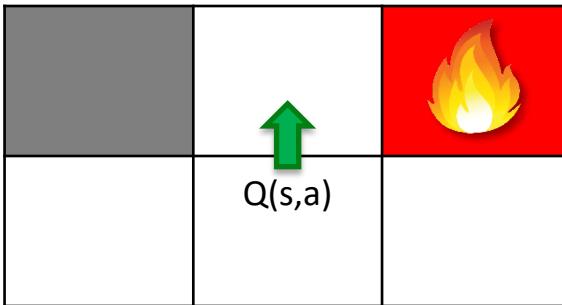


$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$



$$Q(s, a) = R(s, a) + \gamma \sum_{s'} (P(s, a, s') V(s'))$$

Q-Learning Intuição



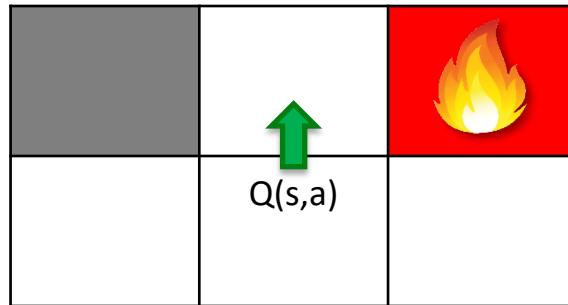
$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} (P(s, a, s') V(s'))$$

Q-Learning Intuição



$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

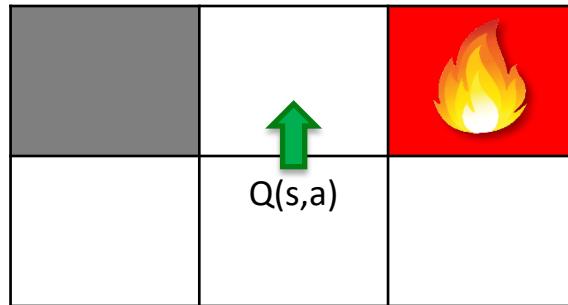


$$Q(s, a) = R(s, a) + \gamma \sum_{s'} \left(P(s, a, s') V(s') \right)$$

Q-Learning Intuição



$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$

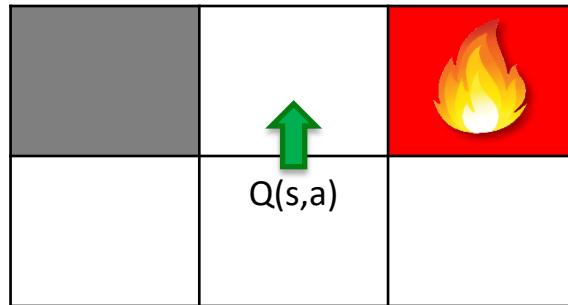


$$Q(s, a) = R(s, a) + \gamma \sum_{s'} \left(P(s, a, s') V(s') \right)$$

Q-Learning Intuição



$$V(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s, a, s') V(s') \right)$$



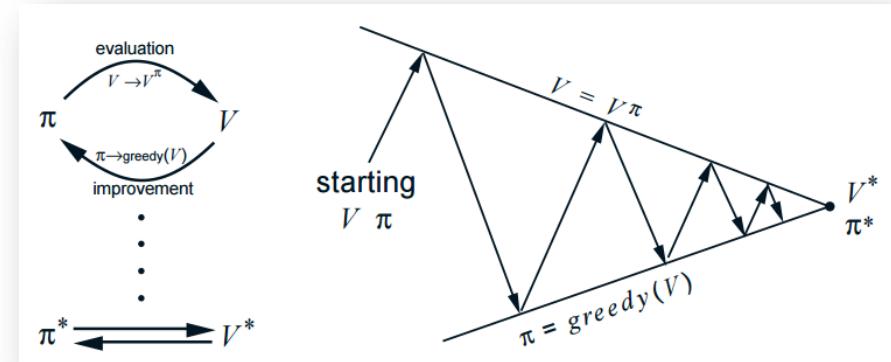
$$Q(s, a) = R(s, a) + \gamma \sum_{s'} \left(P(s, a, s') \max_{a'} Q(s', a') \right)$$

Leitura Adicional

Leitura Adicional:

*Markov Decision Processes:
Concepts and Algorithms*

Martijn van Otterlo (2009)



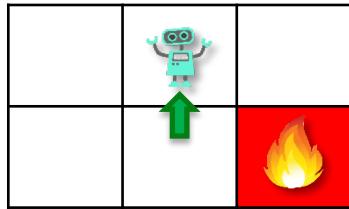
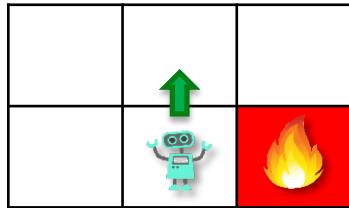
Link:

<https://pdfs.semanticscholar.org/968b/ab782e52faf0f7957ca0f38b9e9078454afe.pdf>

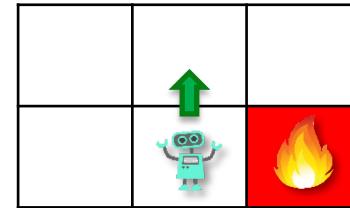
Diferença Temporal

Diferença Temporal

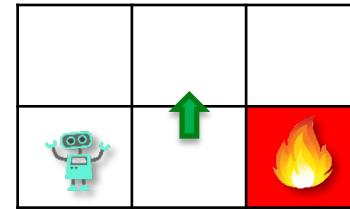
Busca Determinística



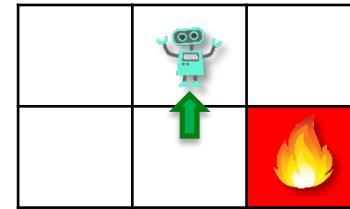
Busca Não-Determinística



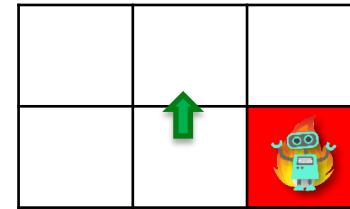
10%



80%



10%



Diferença Temporal

$V=0.81$	$V=0.9$	$V=1$	
$V=0.73$		$V=0.9$	
$V=0.66$	$V=0.73$	$V=0.81$	$V=0.73$

Diferença Temporal

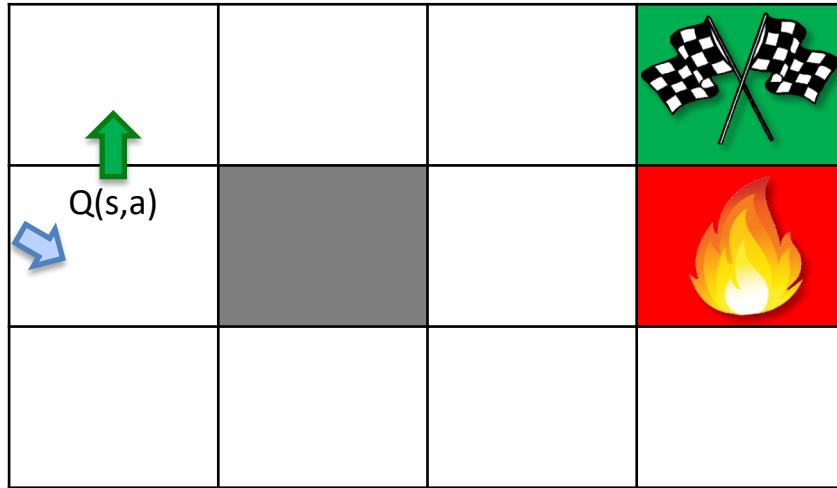
V=0.71	V=0.74	V=0.86	
V=0.63		V=0.39	
V=0.55	V=0.46	V=0.36	V=0.22

Diferença Temporal

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} \left(P(s, a, s') \max_{a'} Q(s', a') \right)$$

$$Q(s, a) = R(s, a) + \gamma \underbrace{\max_{a'} Q(s', a')}$$

Diferença Temporal



Antes:

$$Q(s, a)$$

Depois:

$$R(s, a) + \gamma \max_{a'} Q(s', a')$$

$$TD(a, s) = R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)$$

Diferença Temporal

$$TD(a, s) = R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)$$

$$Q(s, a) = Q(s, a) + \alpha TD(a, s)$$

Diferença Temporal

$$TD(a, s) = R(s, a) + \gamma \max_{a'} Q(s', a') - Q_{t-1}(s, a)$$

$$Q_t(s, a) = Q_{t-1}(s, a) + \alpha TD_t(a, s)$$

$$Q_t(s, a) = Q_{t-1}(s, a) + \alpha \left(R(s, a) + \gamma \max_{a'} Q(s', a') - Q_{t-1}(s, a) \right)$$

Leitura Adicional

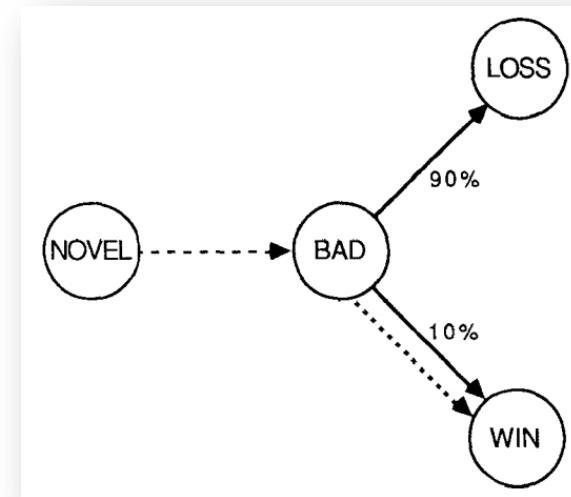
Leitura Adicional:

Learning to Predict by the Methods of Temporal Differences

Richard Sutton (1988)

Link:

<https://link.springer.com/article/10.1007/BF00115009>



Conteúdo

Conteúdo

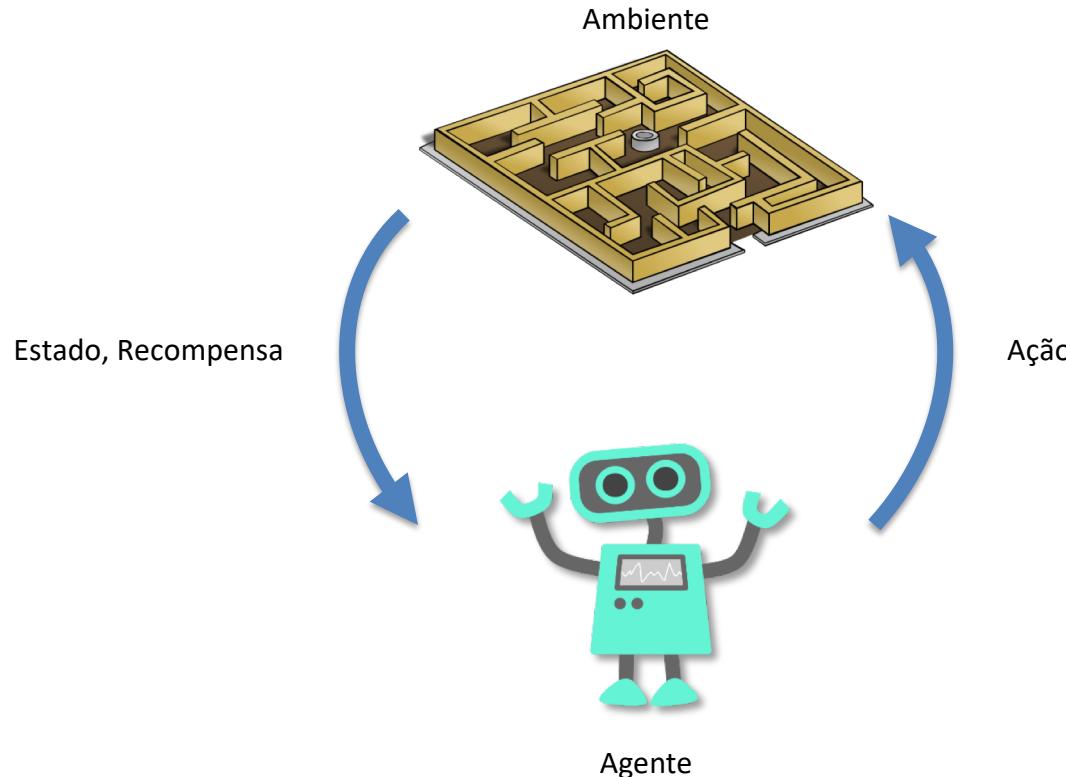
O que você aprenderá nessa seção:

- Intuição Deep Q-Learning (Aprendizagem)
- Intuição Deep Q-Learning (Ação)
- Experiência de Replay
- Políticas de Seleção de Ações

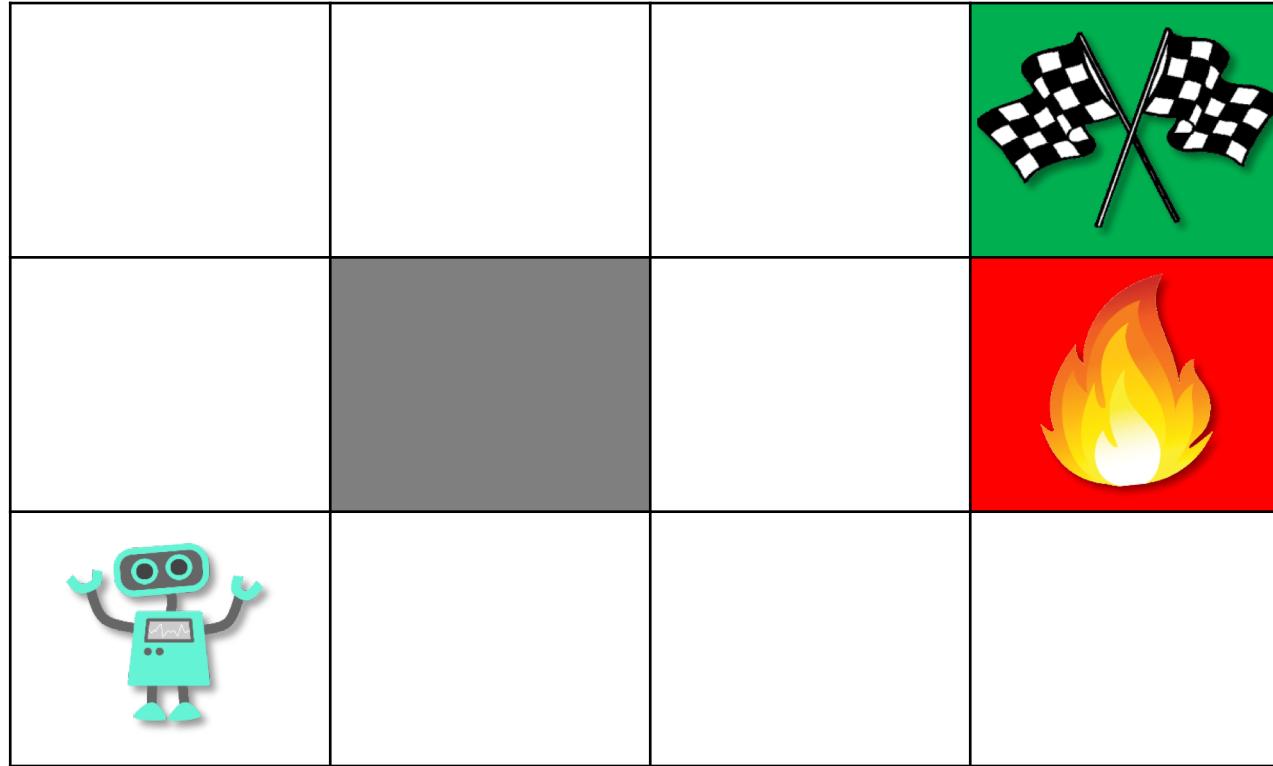
- Anexo 1: Redes Neurais Artificiais

Intuição Deep Q-Learning

Intuição Deep Q-Learning



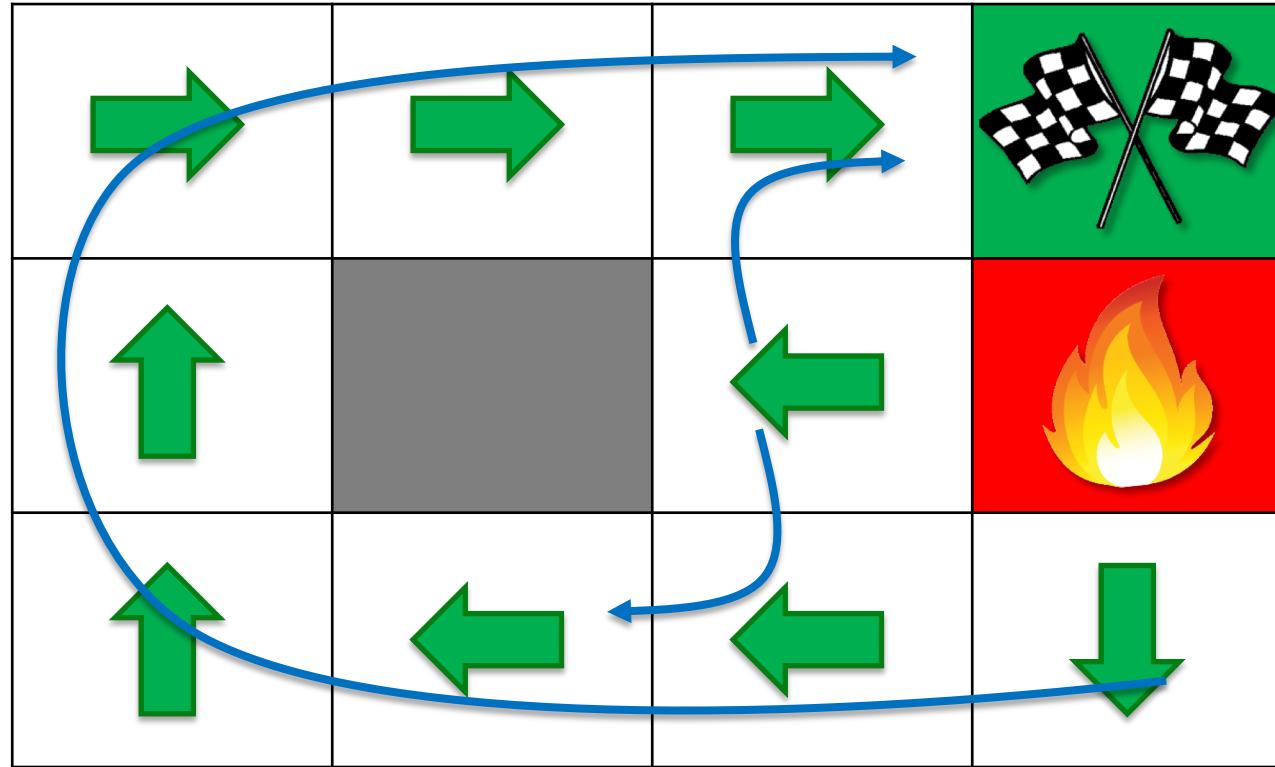
Intuição Deep Q-Learning



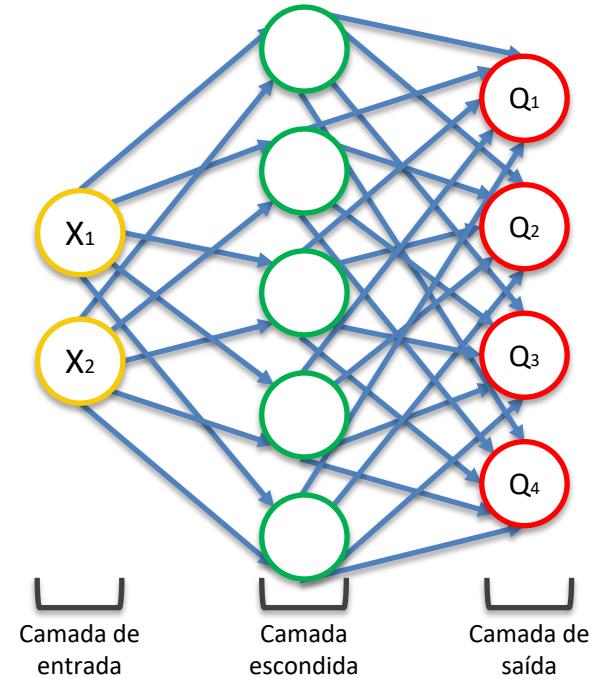
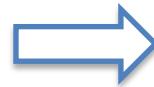
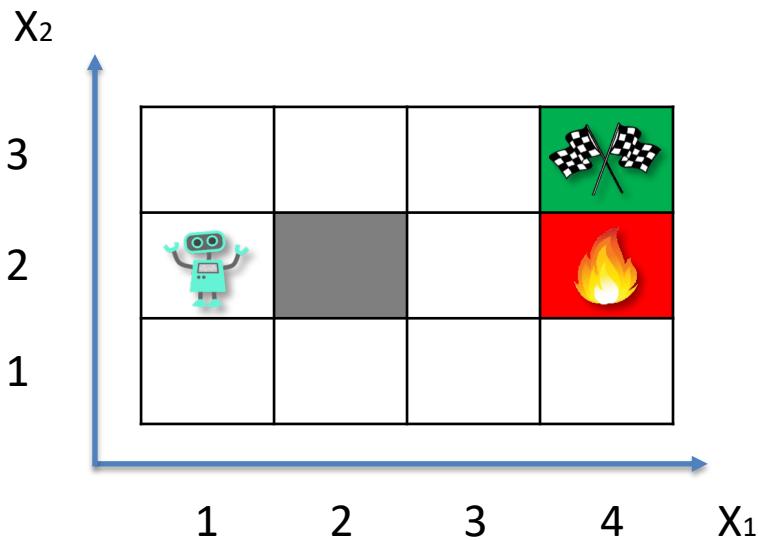
Intuição Deep Q-Learning

$V=0.71$	$V=0.74$	$V=0.86$	
$V=0.63$		$V=0.39$	
$V=0.55$	$V=0.46$	$V=0.36$	$V=0.22$

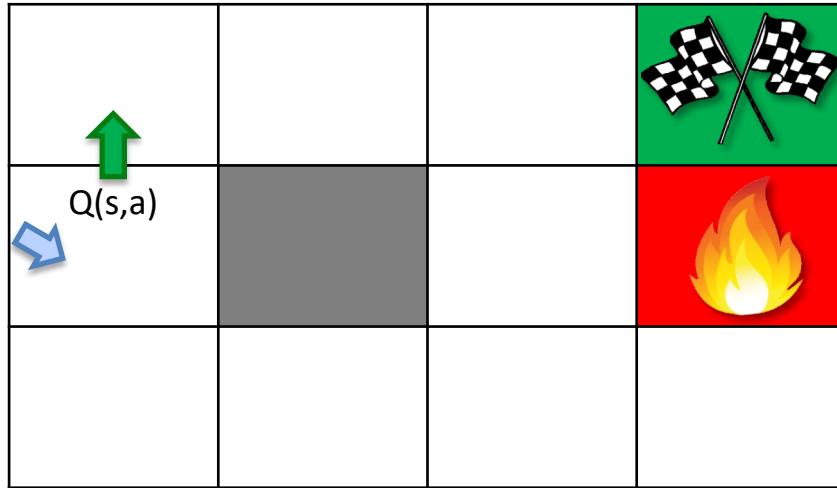
Intuição Deep Q-Learning



Intuição Deep Q-Learning



Intuição Deep Q-Learning



Antes:

$$Q(s, a)$$

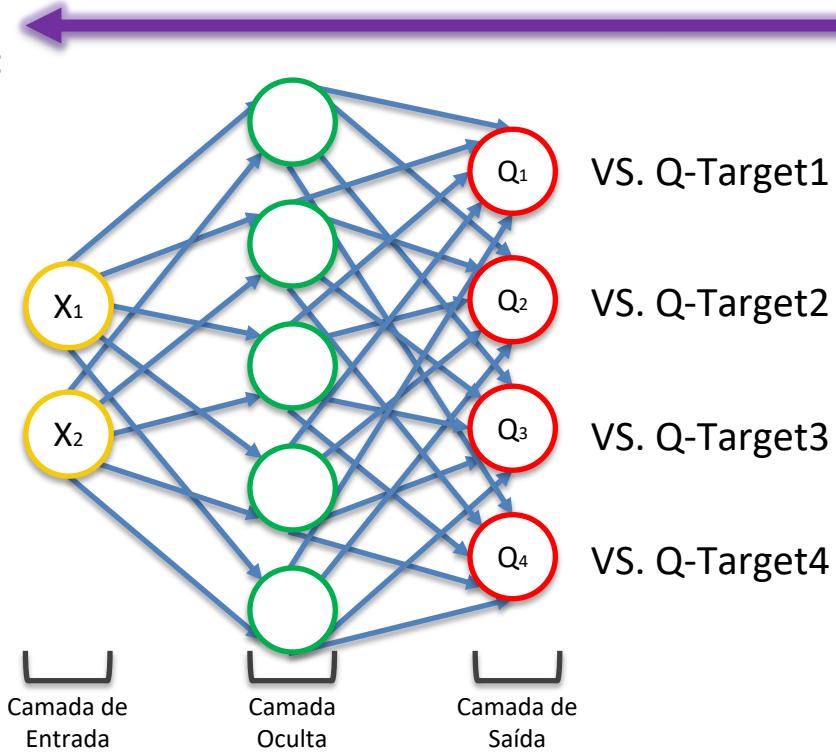
Depois:

$$R(s, a) + \gamma \max_{a'} Q(s', a')$$

$$TD(a, s) = R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)$$

Intuição Deep Q-Learning

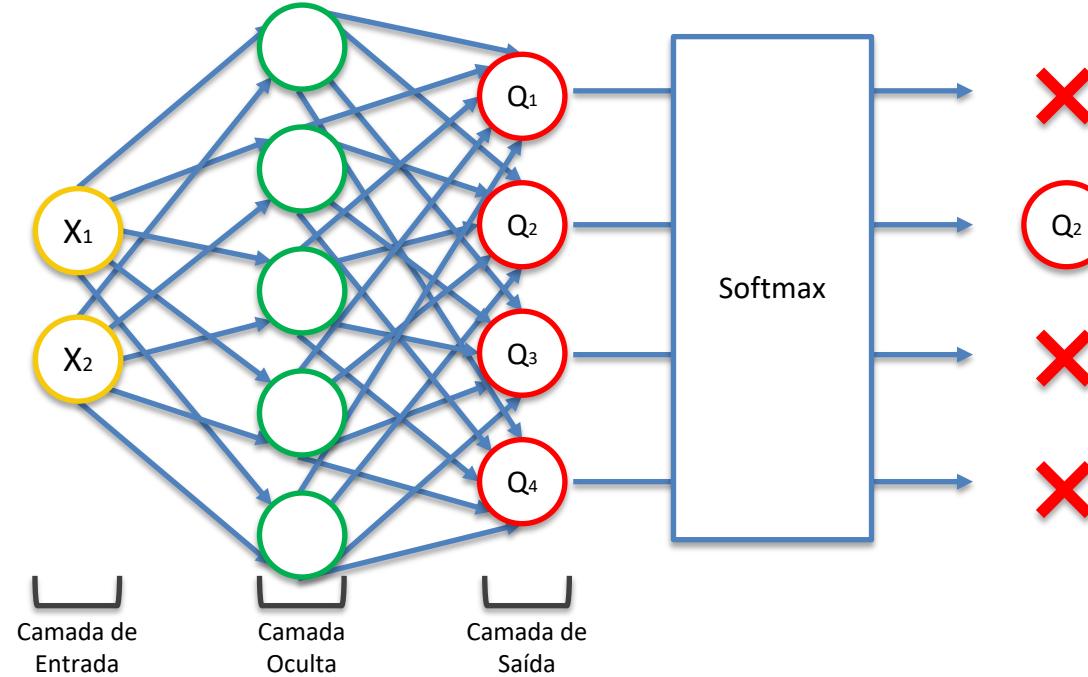
Aprendizagem:



$$L = \sum (Q\text{-}Target - Q)^2$$

Intuição Deep Q-Learning

Ações:



Intuição Deep Q-Learning

Aprendizagem:

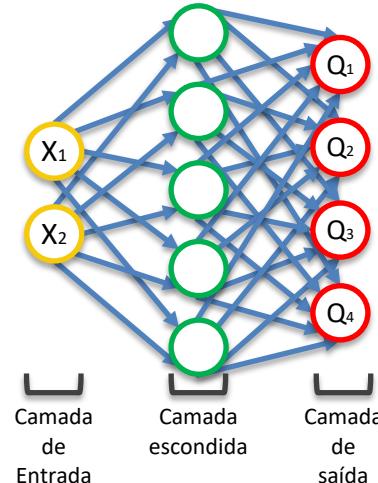
VS. Q-Target1

VS. Q-Target2

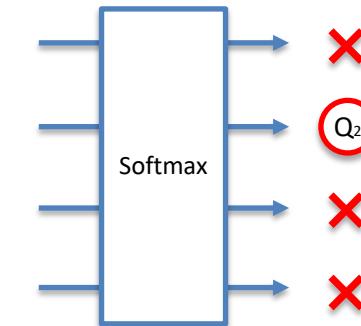
VS. Q-Target3

VS. Q-Target4

$$L = \sum (Q\text{-}Target - Q)^2$$



Ações:

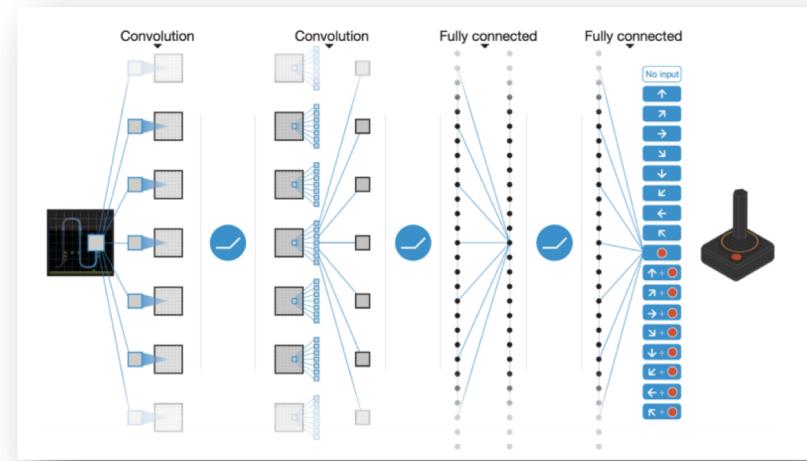


Leitura Adicional

Leitura Adicional:

Simple Reinforcement Learning with Tensorflow (Part 4)

Arthur Juliani (2016)



Link:

<https://medium.com/@awjuliani/simple-reinforcement-learning-with-tensorflow-part-4-deep-q-networks-and-beyond-8438a3e2b8df>

Experiência de Replay

Experiência de Replay

Aprendizagem:

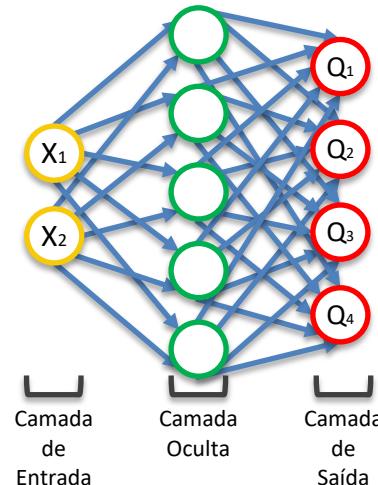
VS. Q-Target1

VS. Q-Target2

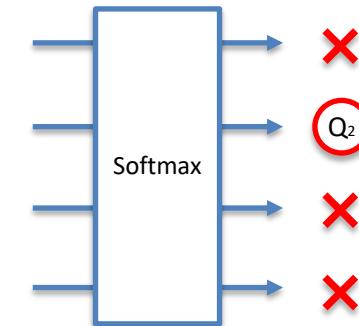
VS. Q-Target3

VS. Q-Target4

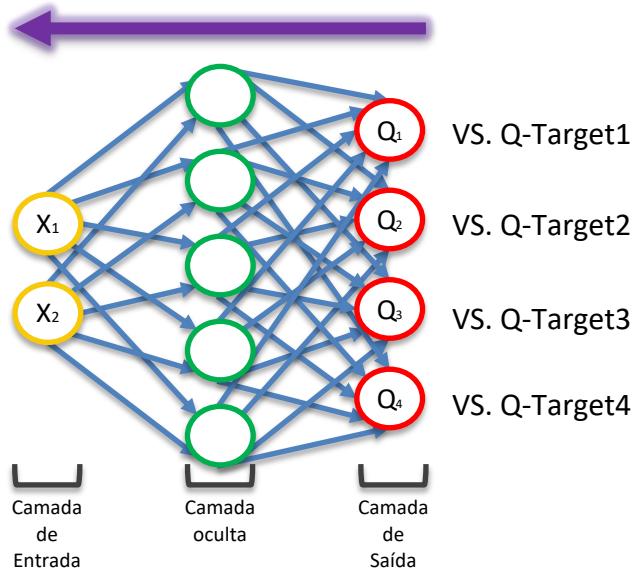
$$L = \sum (Q\text{-}Target - Q)^2$$



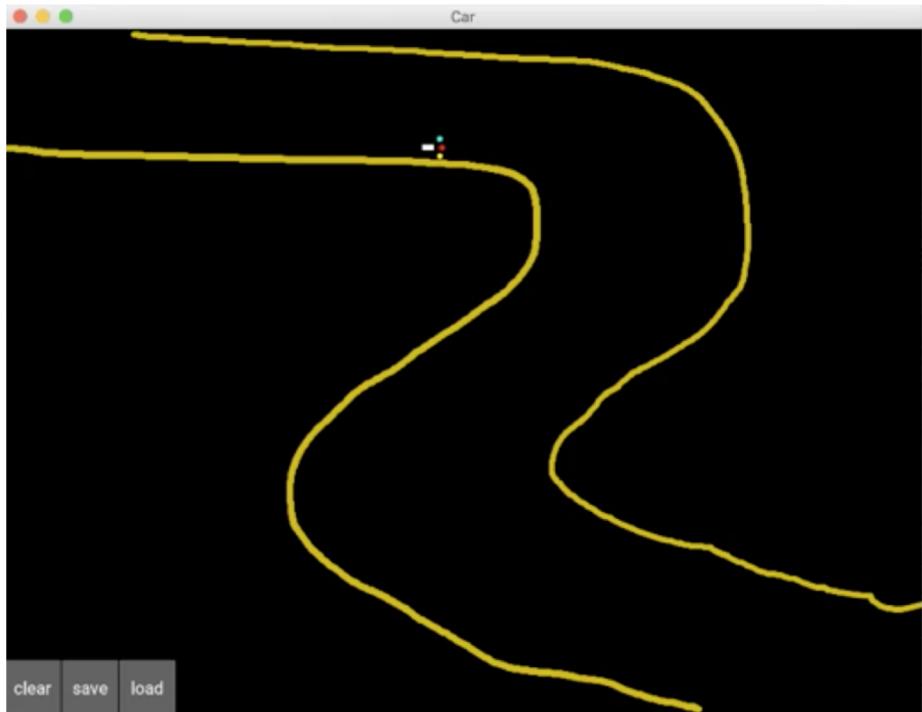
Ação:



Experiência de Replay



$$L = \sum (Q\text{-Target} - Q)^2$$

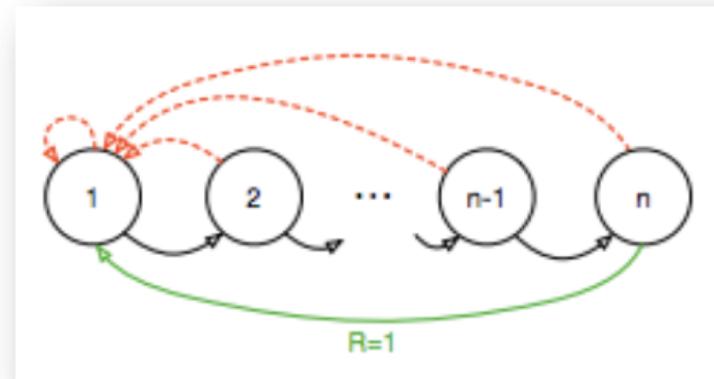


Leitura Adicional

Leitura Adicional:

Prioritized Experience Replay

Tom Schaul et al.,
Google DeepMind (2016)



Link:

<https://arxiv.org/pdf/1511.05952.pdf>

Política de Seleção de Ações

Política de Seleção de Ações

Aprendizagem:

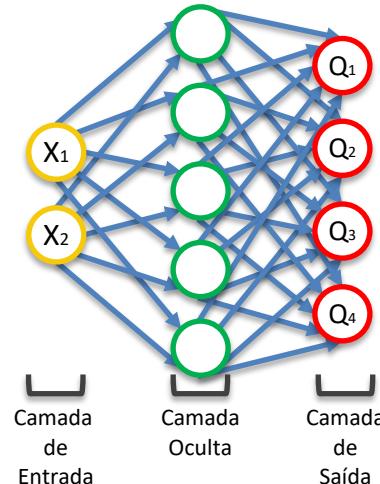
VS. Q-Target1

VS. Q-Target2

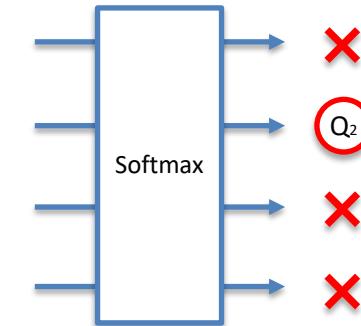
VS. Q-Target3

VS. Q-Target4

$$L = \sum (Q\text{-}Target - Q)^2$$

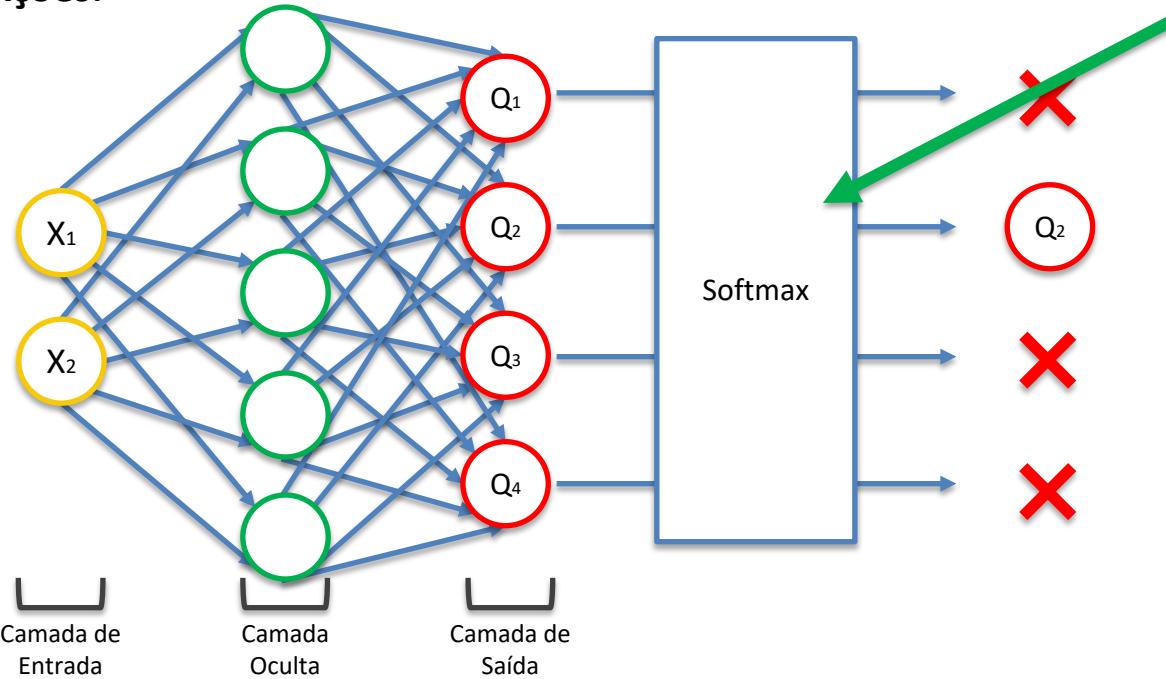


Ações:



Política de Seleção de Ações

Ações:

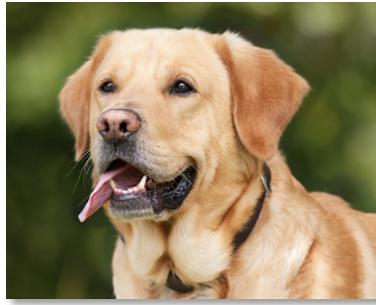


Seleção de Ação:

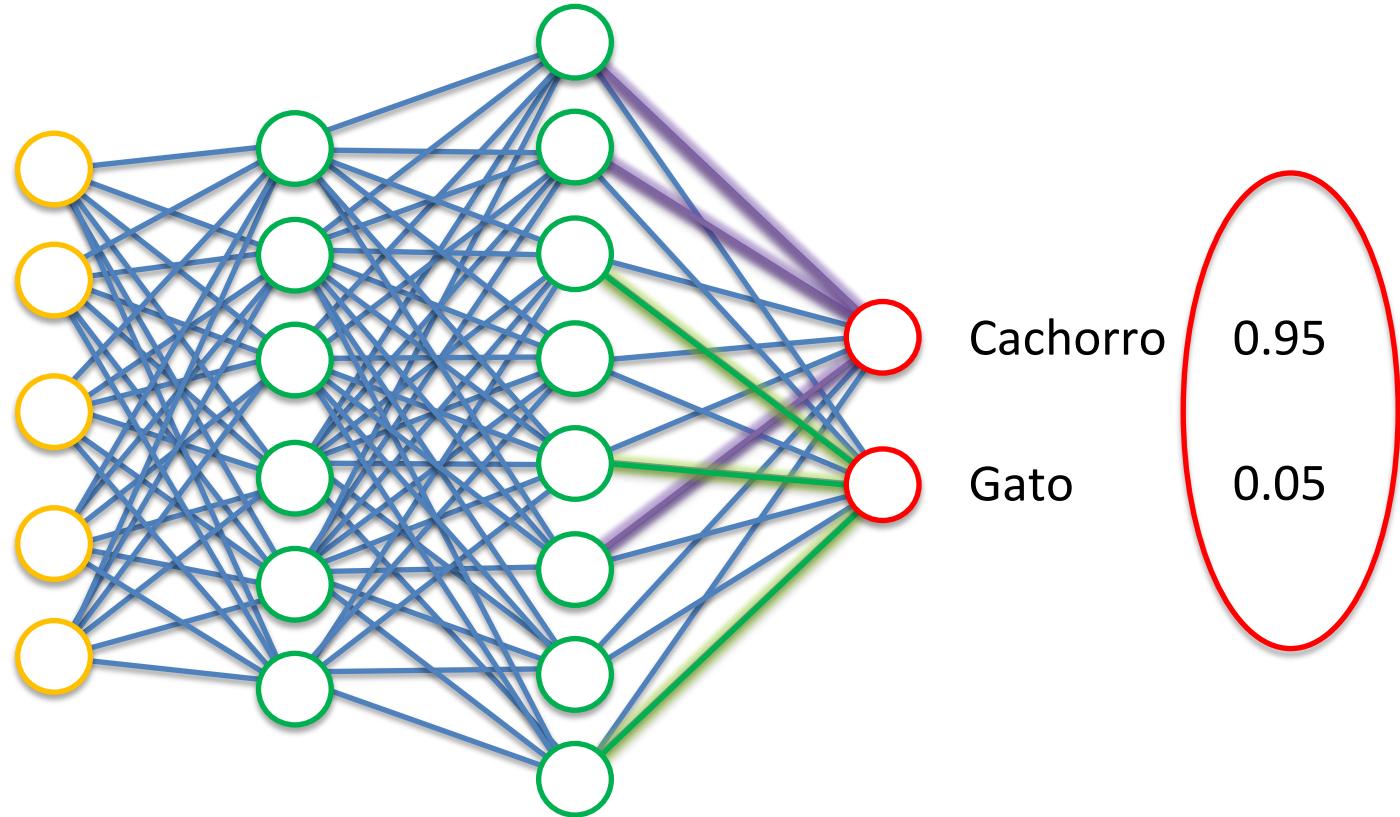
- ϵ -greedy
- ϵ -soft ($1-\epsilon$)
- Softmax

Exploration
vs
Exploitation

Política de Seleção de Ações



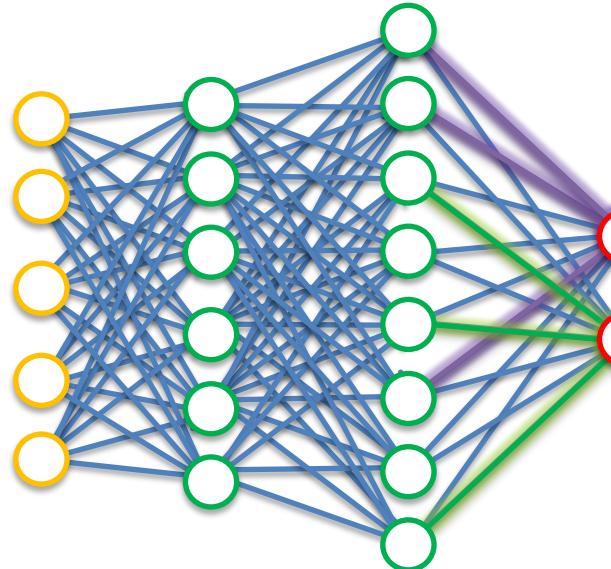
Flattening



Política de Seleção de Ações



.....
Flattening

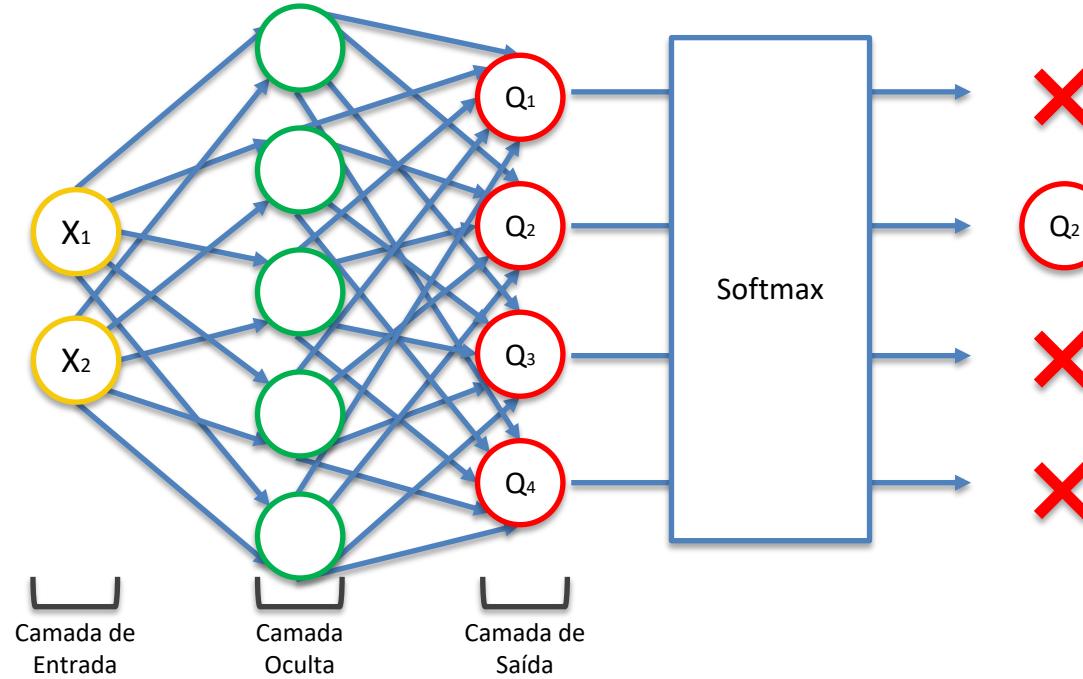


$$f_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

Cachorro $\longrightarrow z_1 \longrightarrow 0.95$
Gato $\longrightarrow z_2 \longrightarrow 0.05$

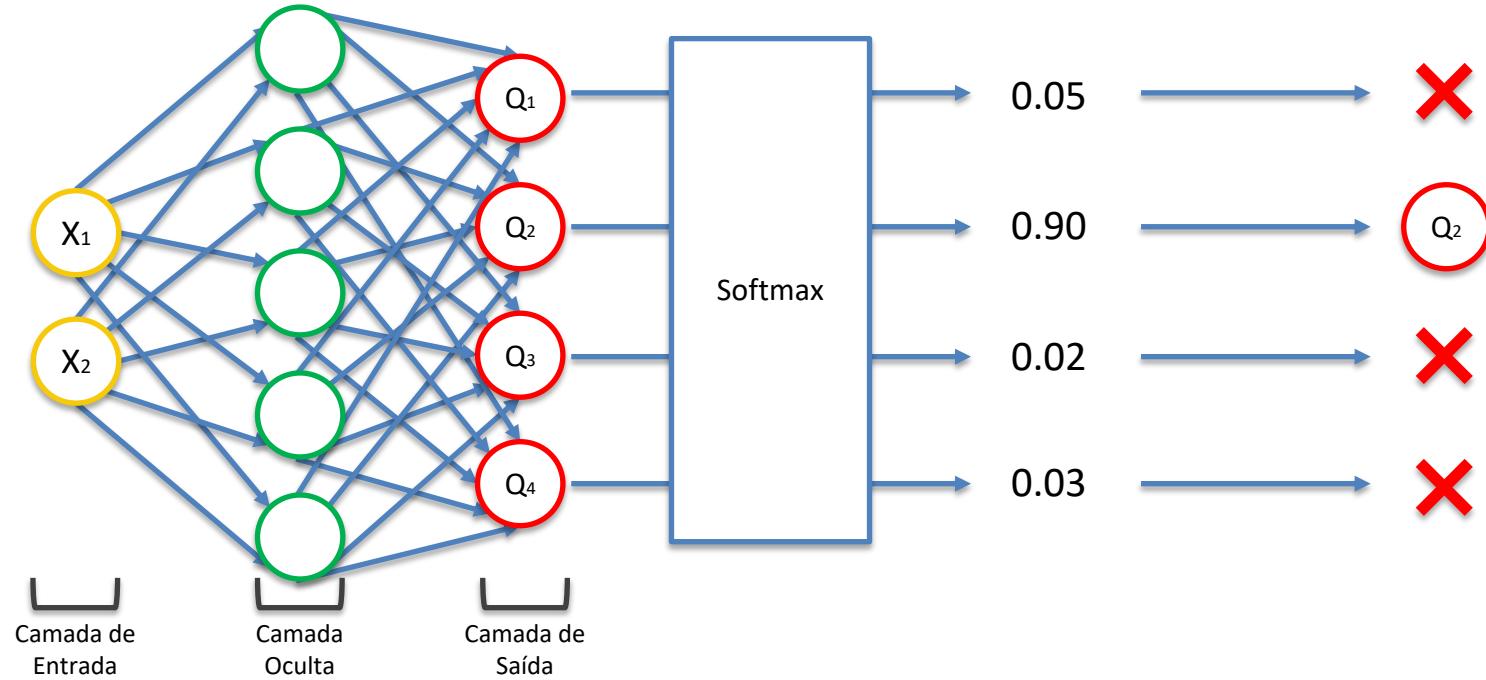
Política de Seleção de Ações

Ações:



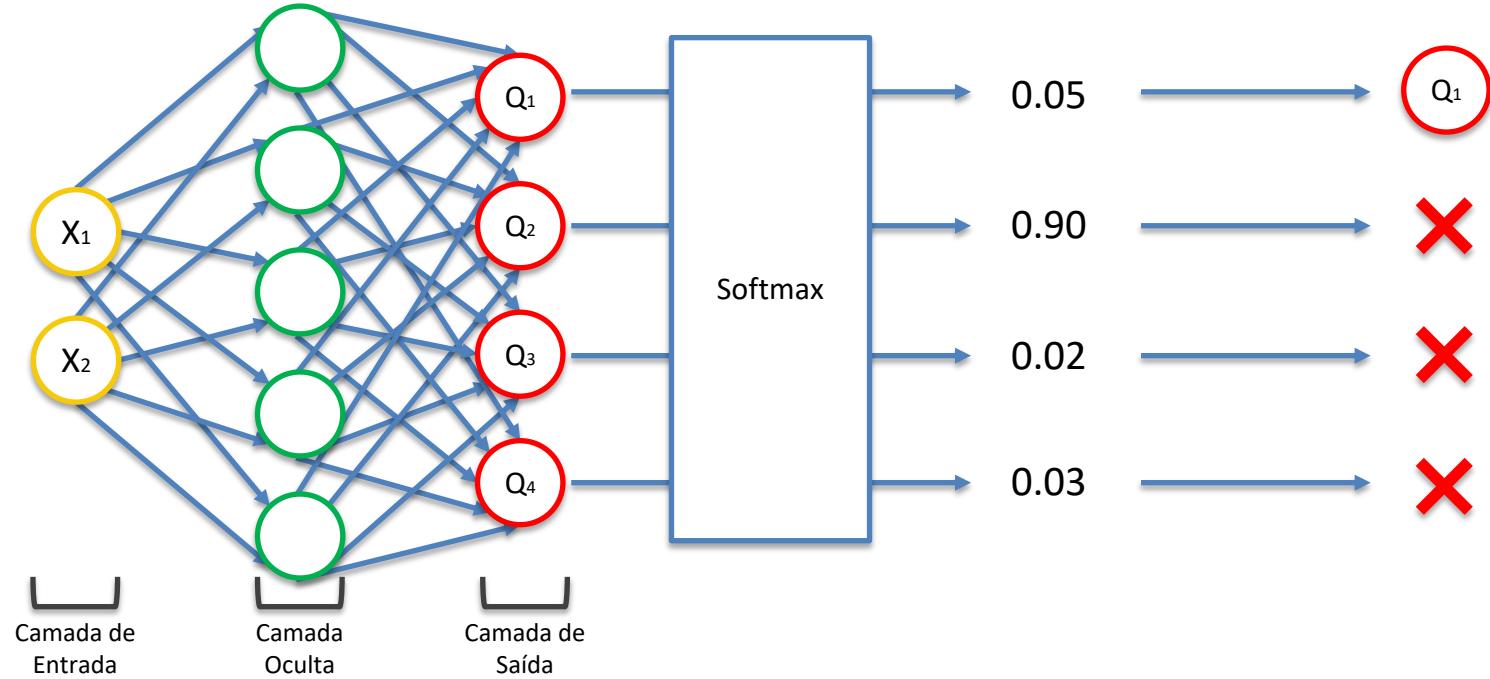
Política de Seleção de Ações

Ações:



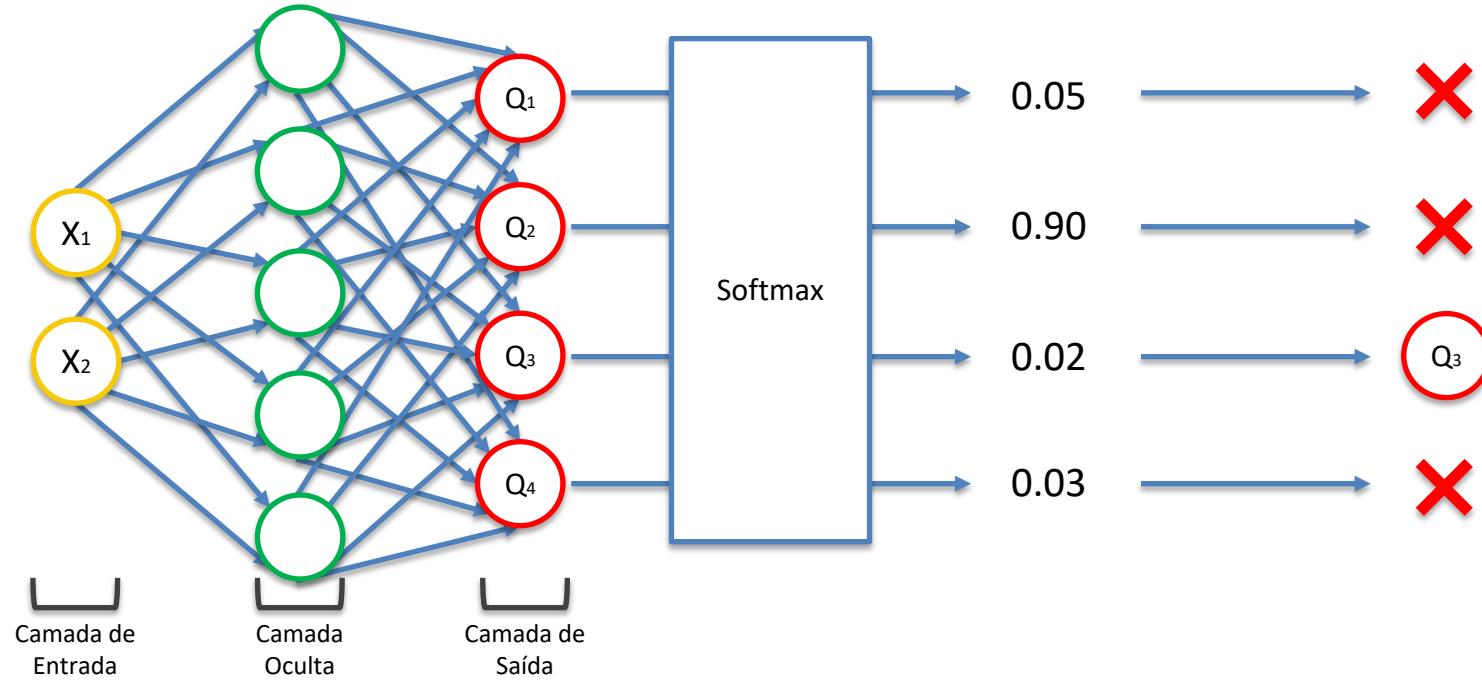
Política de Seleção de Ações

Ações:



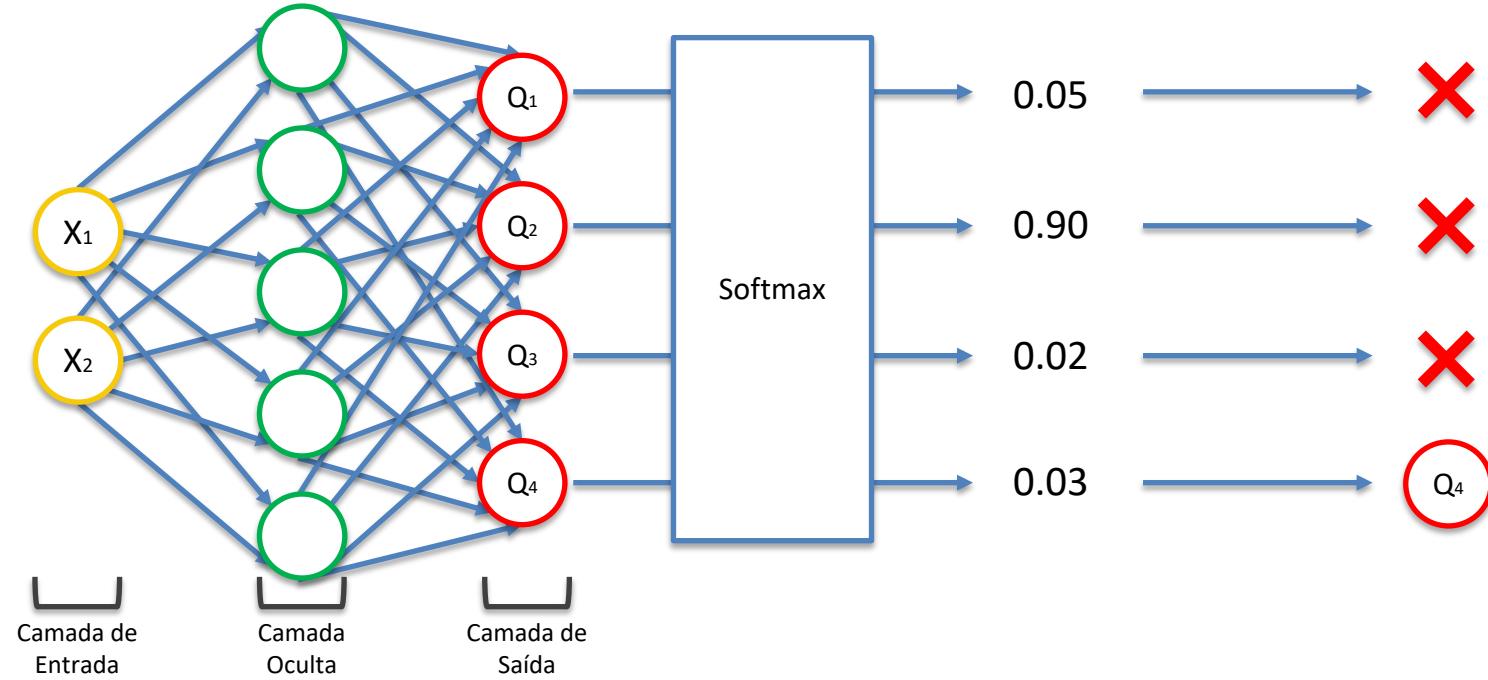
Política de Seleção de Ações

Ações:



Política de Seleção de Ações

Ações:



Leitura Adicional

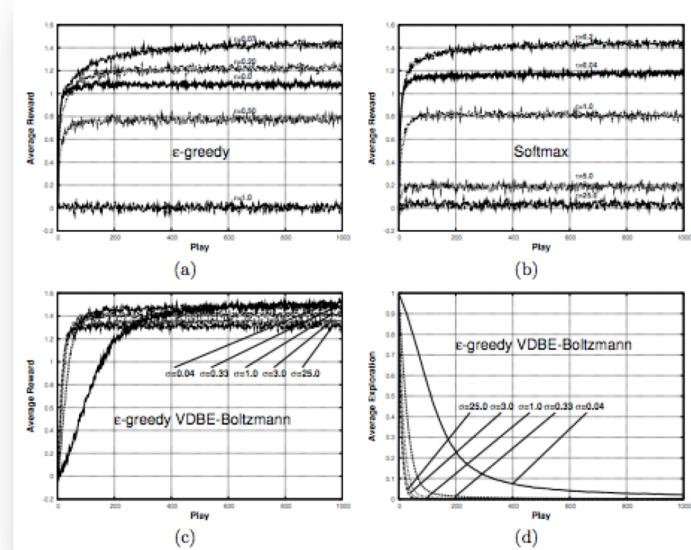
Leitura Adicional:

Adaptive ϵ -greedy Exploration in Reinforcement Learning Based on Value Differences

Michel Tokic (2010)

Link:

<http://tokic.com/www/tokicm/publikationen/papers/AdaptiveEpsilonGreedyExploration.pdf>



Conteúdo

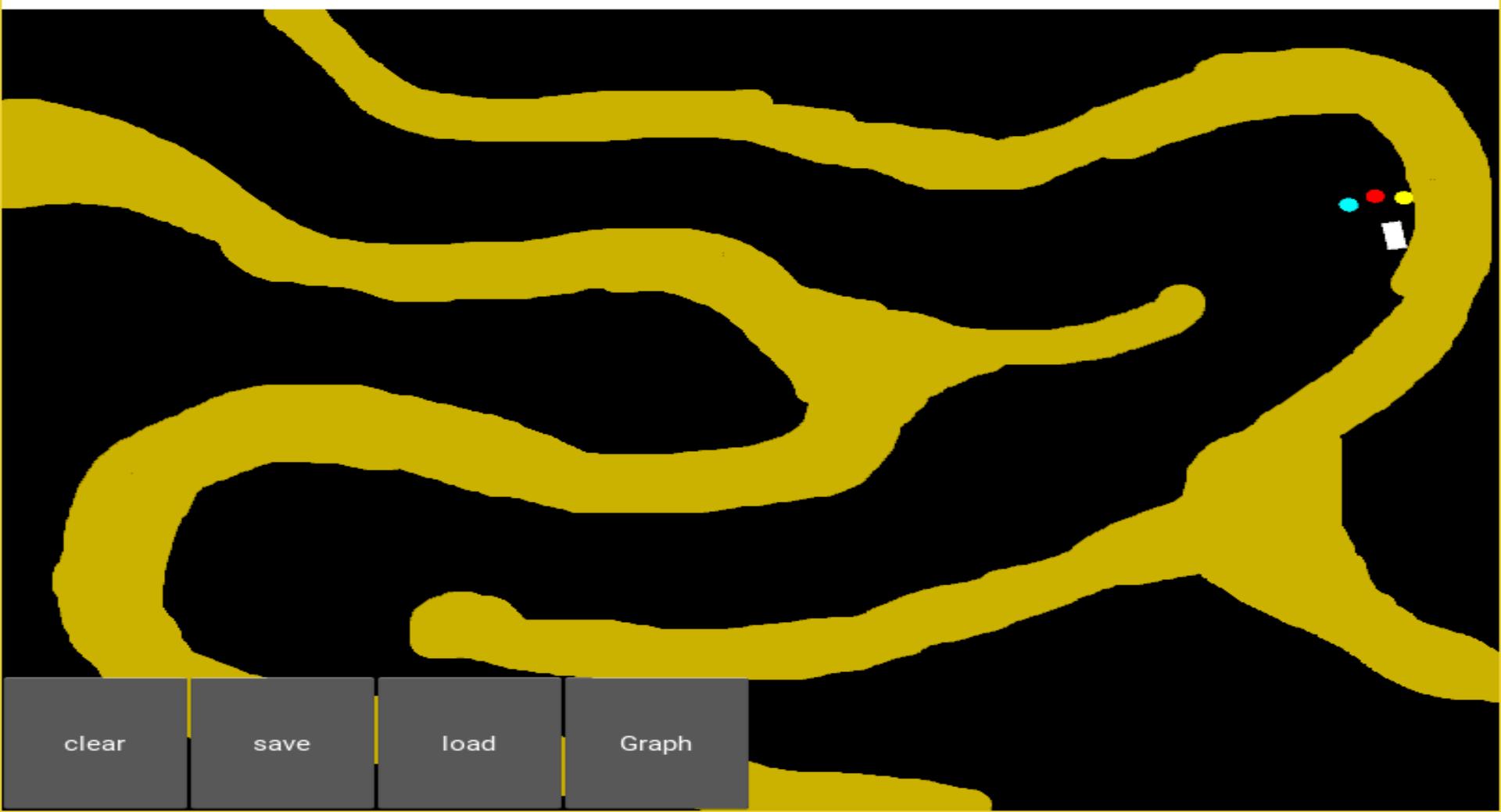
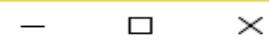
Conteúdo

O que você aprenderá nessa seção:

- Construção de um carro autônomo!
- Ambiente contendo o mapa
- IA para controlar o carro no mapa com Deep Q-Learning (PyTorch)
- Configuração do Ambiente



Car

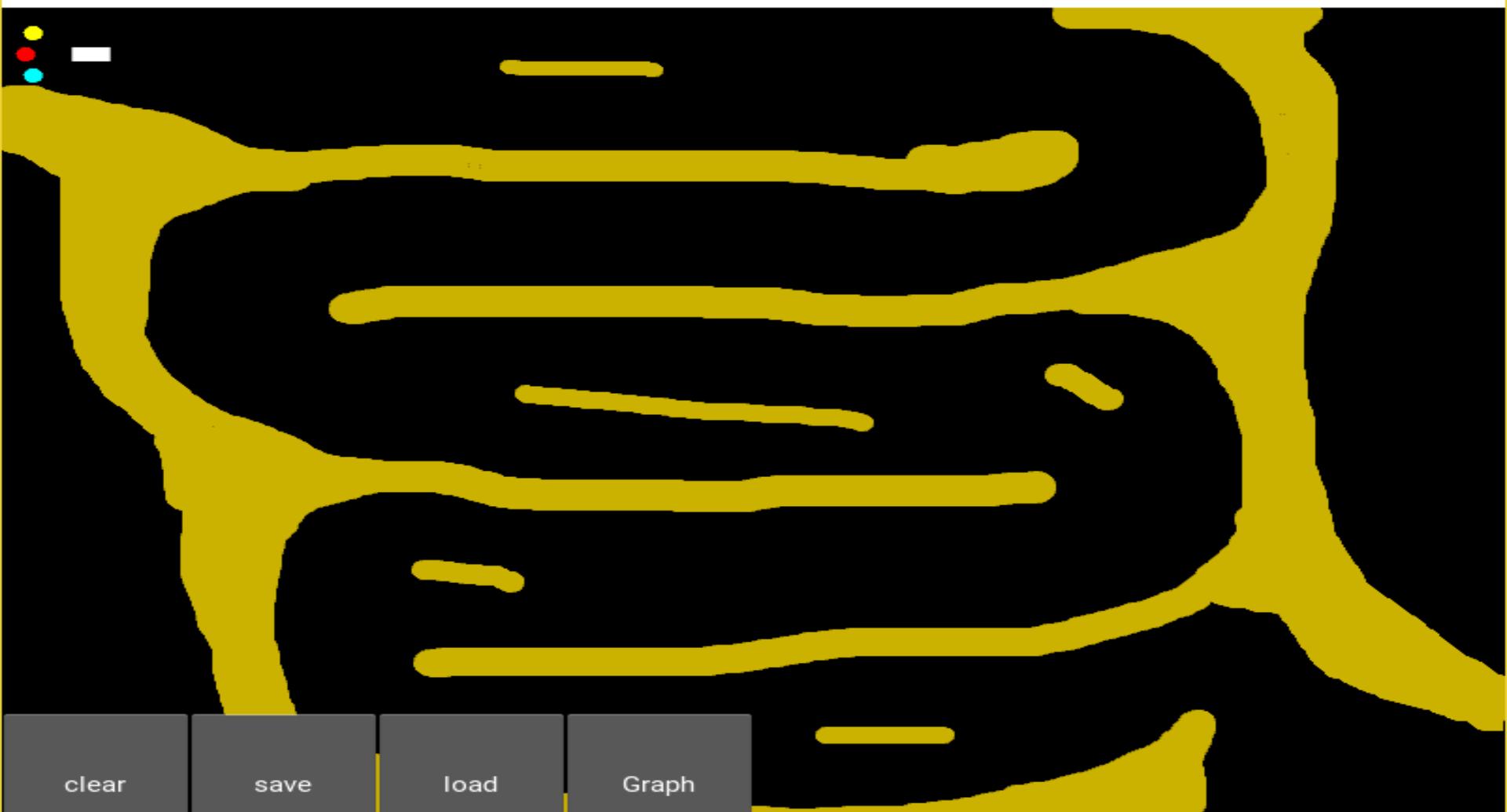


clear

save

load

Graph



clear

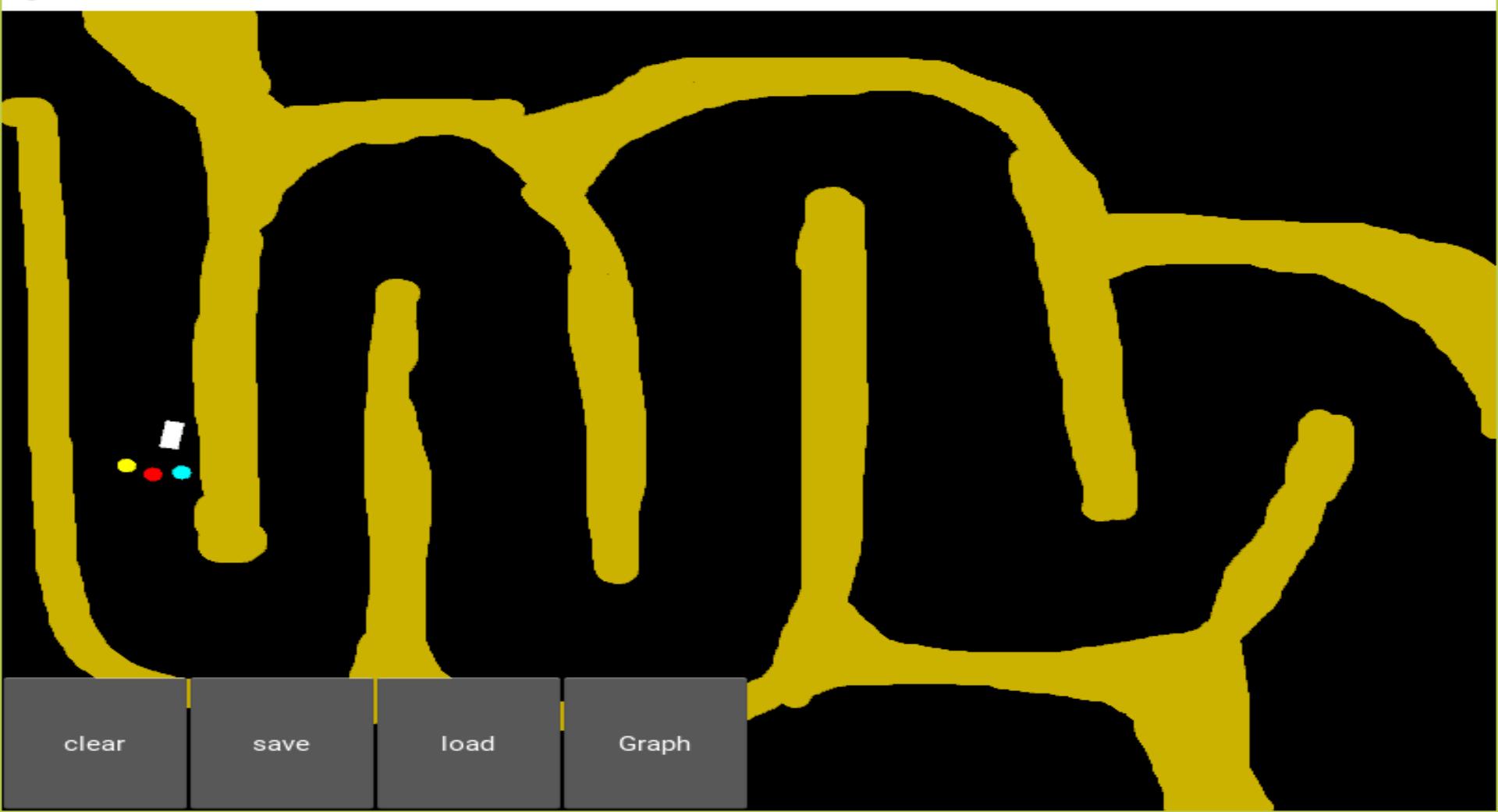
save

load

Graph



Car



clear

save

load

Graph

