



Aprendizado por Reforço

AULA - 1

Conceitos e Modelagem

Anteriormente...



- Otimizar: Encontrar a melhor solução para um problema
 - Maximizar ou minimizar uma função é encontrar os melhores valores para suas variáveis
- Aprendizado de máquina: Reconhecimento de padrões e fazer associações sem programação explícita
- Redes Neurais Artificiais: Conjuntos de neurônios paralelos e sequenciais cujas variáveis são alteradas durante o treinamento.
- **Imitation Learning**: Aprender comportamentos a partir de exemplos de experts
 - Problemas sequenciais precisam de representatividade nos dados, especialmente para correção de caminhos errados.

O que é aprender por REFORÇO?

- Psicologia
 - Behaviorismo
- Adestração Animal



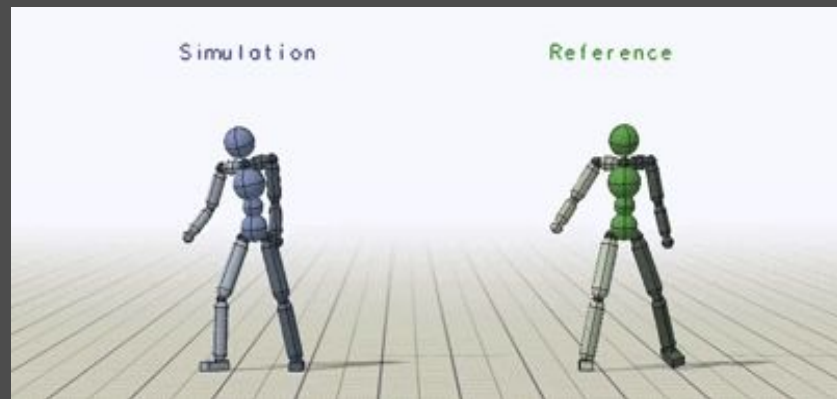
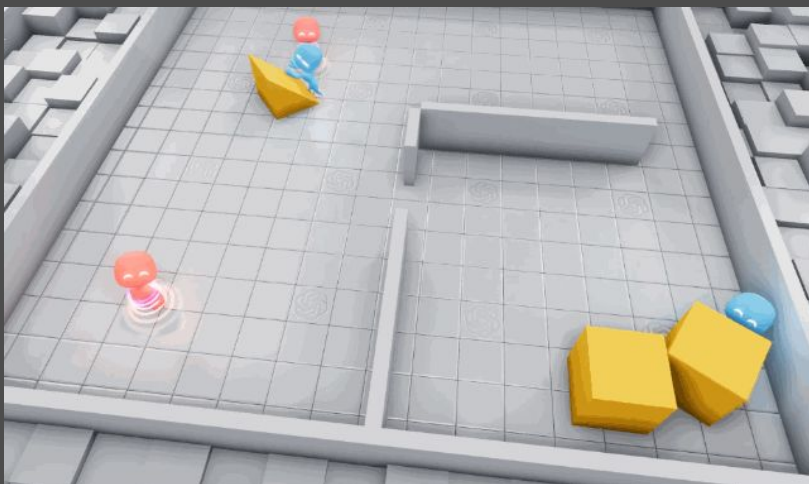
Mapeamento Estados e Ações

Estado/Situação	Melhor Ação
Ouvir "Senta!"	Sentar
Ouvir "Pega!"	Atacar Alvo
Ouvir "Dá a pata"	Erguer pata dianteira



“In interactive problems it is often impractical to obtain examples of desired behavior that are both correct and representative of all the situations in which the agent has to act.”

Outros Exemplos

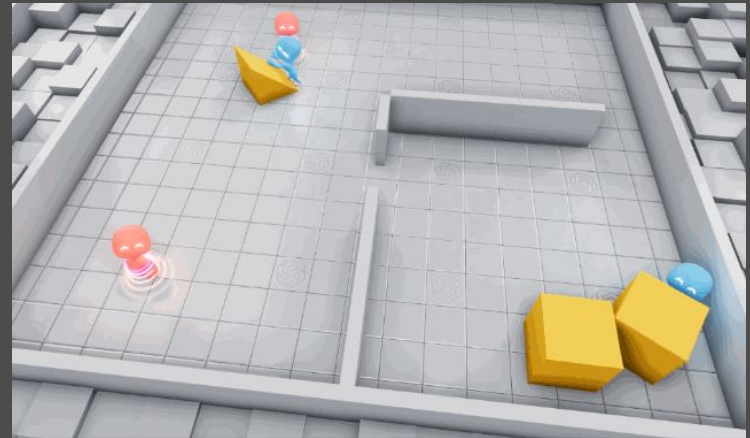


Relação Agente e Ambiente



Ambiente

- Elemento com o qual o agente irá interagir.
 - Engloba tudo que não é o Agente
 - Inclusive possíveis outros agentes



Observação/Estado

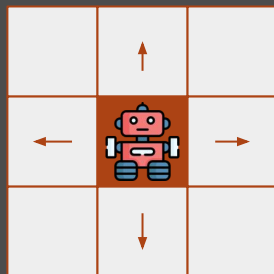
- Descrição da situação atual do ambiente
 - Pode ser parcial
 - Leitura de um sensor
 - Imagem
- Conjunto de todos os estados possíveis:
- Espaço de Estados



Observação deve permitir a inferência da recompensa

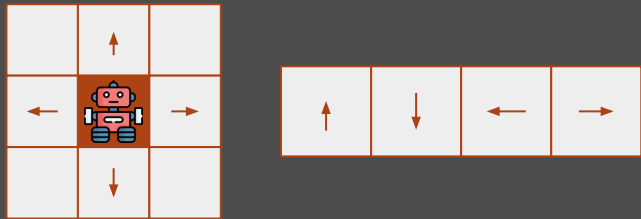
Ações

- Forma do Agente interagir com o Ambiente
- Conjunto de todas as ações possíveis:
- Espaço de Ações



Espaço de Ações

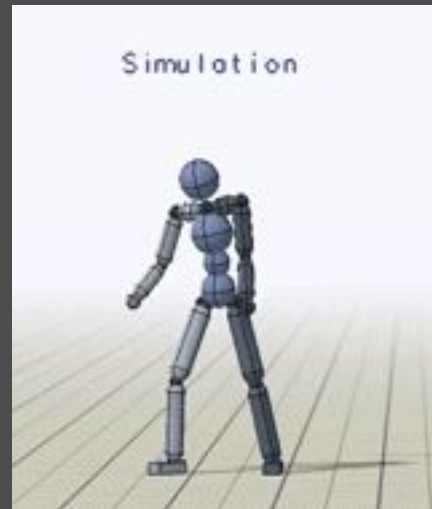
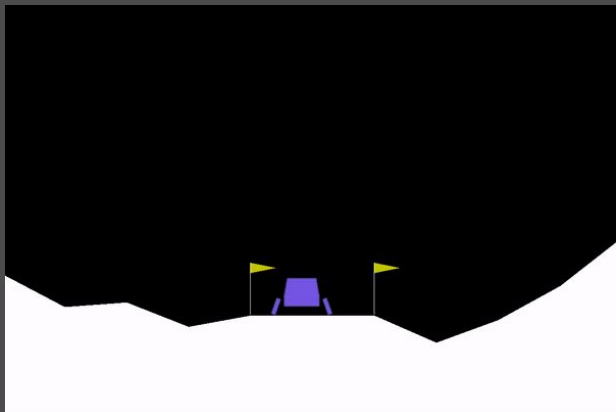
- Escolher uma dentre N ações: espaço unidimensional



- Escolher mais de uma ação simultânea: uma dimensão para cada ação que será escolhida

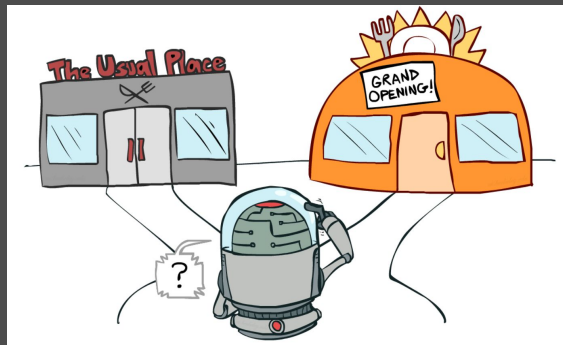


Qual o Espaço de Ações para estes problemas?



Exploration vs Exploitation

- **Exploration:** Encontrar mais informação sobre o ambiente, é preciso testar ações novas.
- **Exploitation:** Utilizar da informação adquirida previamente para maximizar a recompensa, é preciso maximizar a recompensa.
- Quanto mais opções de **ações**, mais difícil fica de explorar soluções.
- Quanto mais **situações** diferentes, também fica mais difícil de explorar o ambiente.



A complexidade do problema vai ser definida pela soma das dimensões do espaço de estados e do espaço de ações.



A questão não é apenas quantas opções eu tenho, mas sim em quantas situações eu terei que escolher dentre essas opções

Sinal de Recompensa

- Dizer o quão bom é o estado do ambiente.
- Imediatamente após a entrada no estado.
- “Atrelado” à última ação tomada.
- Diz o quão boa foi aquela ação específica?
 - Não necessariamente

s, a, r, s'

Transição



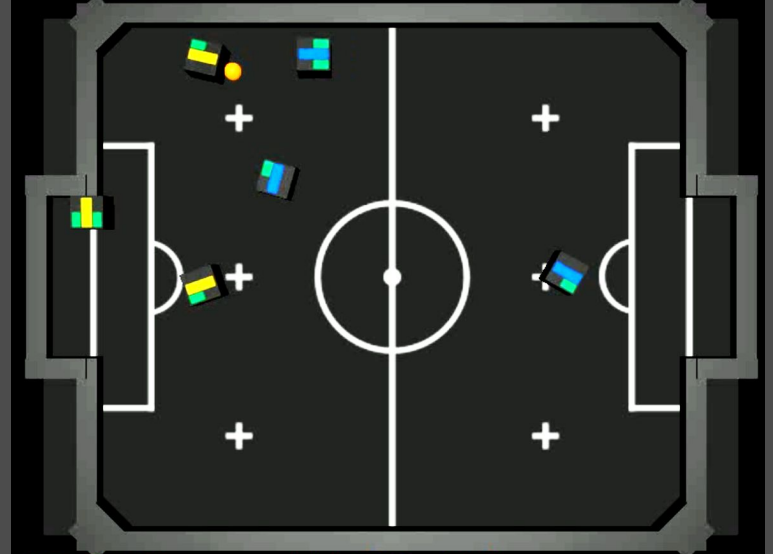
Função de Recompensa

Qual seria uma boa função de recompensa para um jogo de xadrez?



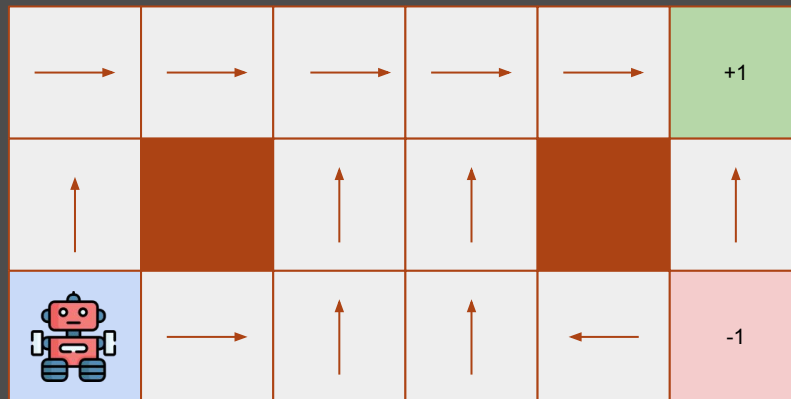
Função de Recompensa

Qual seria uma boa função de recompensa para um jogo de futebol de robôs?



Política

- Comportamento do Agente
- COMO ele mapeia estados a ações
- Representa uma solução para o problema
- O que muda durante o treino



			-0.5	+1

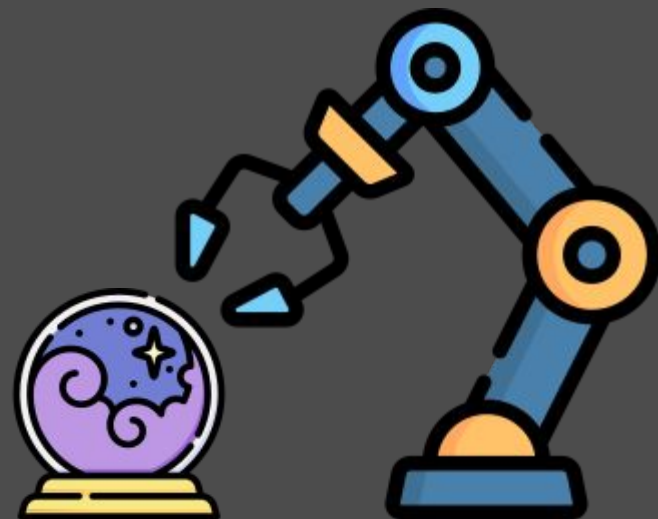
Recompensas

0.22	0.25	0.29	0.31	0.9	1.0
0.2		0.25	0.29		0.9
	0.2	0.22	-0.25	-0.9	-1.0

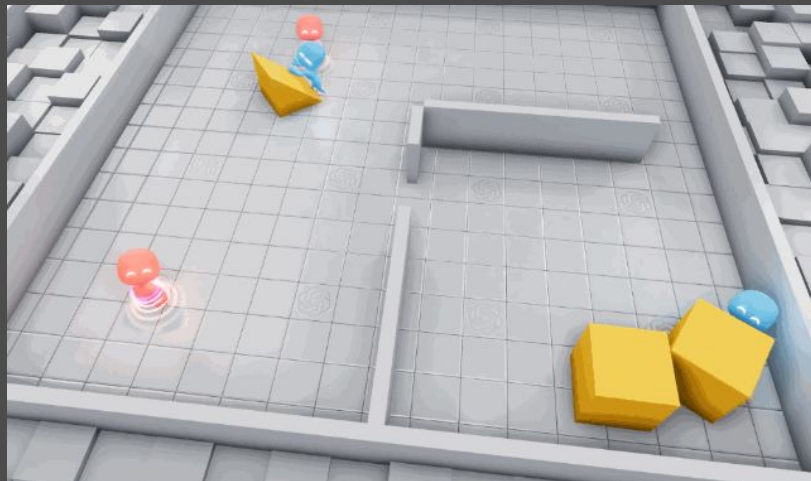
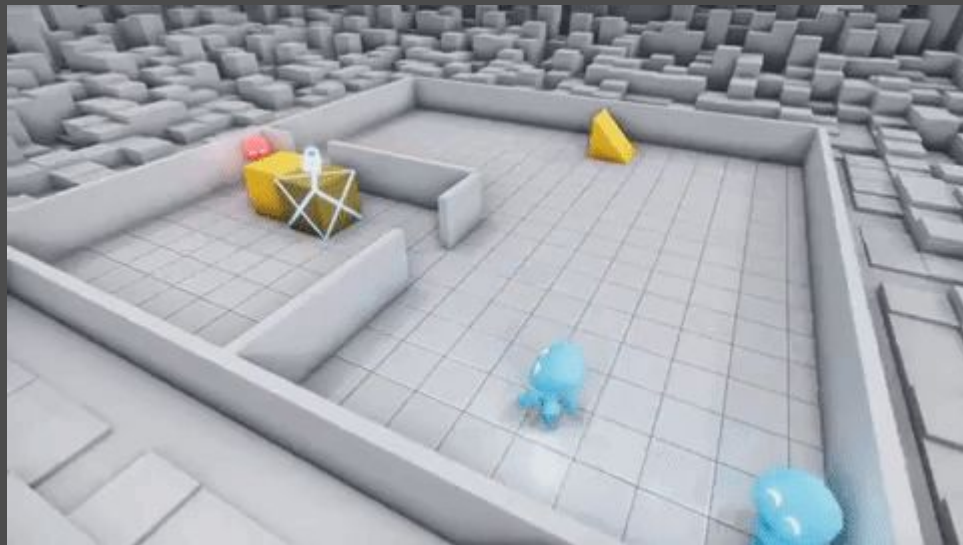
Valores

Modelo (opcional)

- Tenta prever o próximo estado
- “Modelo” do ambiente
- Usados para planejamento
- Auxilia a política
- *Model-based*
- *Model-free*



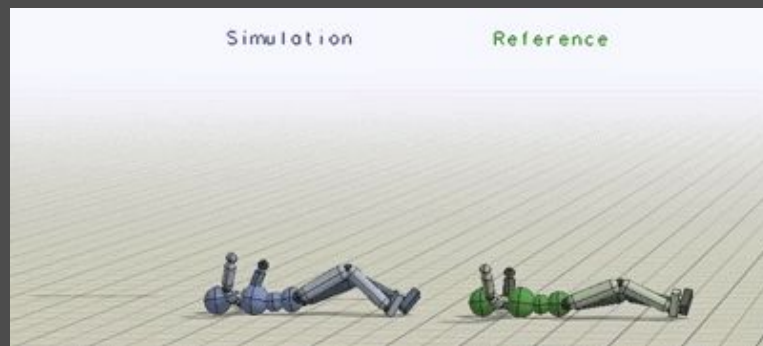
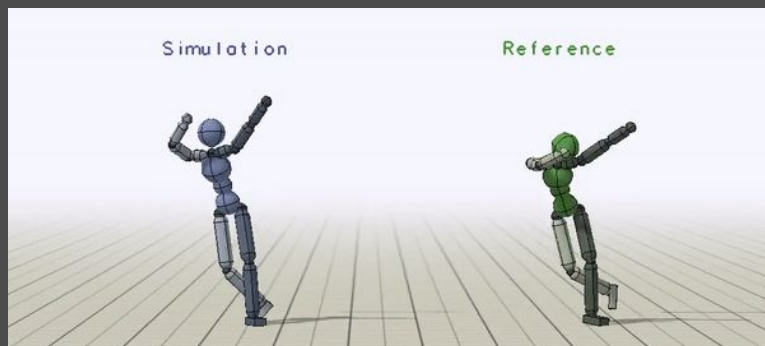
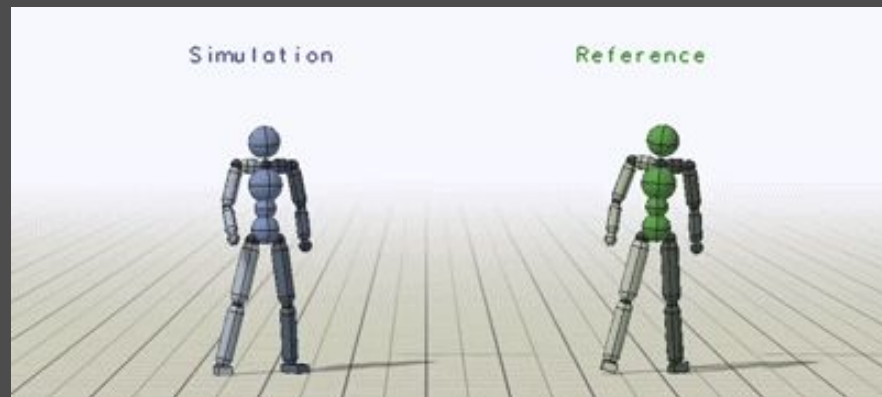
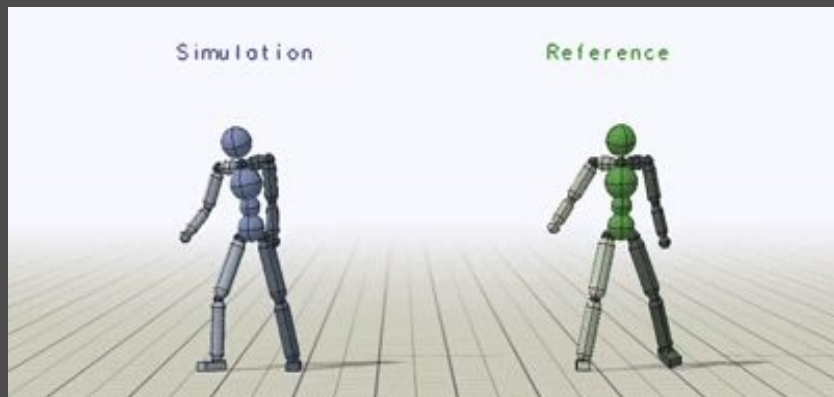
Modelagens:



Modelagens:



Modelagens:





É só...
Por enquanto