

# Instituto Tecnológico de Aeronáutica

---

## **CT-213**

Laboratório 11 - Aprendizado por Reforço Livre de Modelo

---

Aluno: Pedro Elardenberg Sousa e Souza

Professor: Marcos Ricardo Omena de Albuquerque Maximo

21 de maio de 2024

# Conteúdo

<b>1</b>	<b>Resumo</b>	<b>1</b>
<b>2</b>	<b>Introdução</b>	<b>2</b>
2.1	Aprendizado <i>On</i> e <i>Off-Policy</i> . . . . .	2
2.2	Política $\epsilon$ -greedy . . . . .	2
2.3	Sarsa . . . . .	2
2.4	<i>Q-Learning</i> . . . . .	2
<b>3</b>	<b>Análise dos Resultados</b>	<b>3</b>
3.1	Implementação dos algoritmos de RL . . . . .	3
3.2	Aprendizado da política do robô seguidor de linha . . . . .	4
<b>4</b>	<b>Conclusão</b>	<b>6</b>

# 1 Resumo

Neste Laboratório, foram implementados algoritmos de Aprendizado por Reforço (RL) Livre de Modelo, a saber *Sarsa* e *Q-Learning*, e com eles resolvido o problema de um robô seguidor de linha.

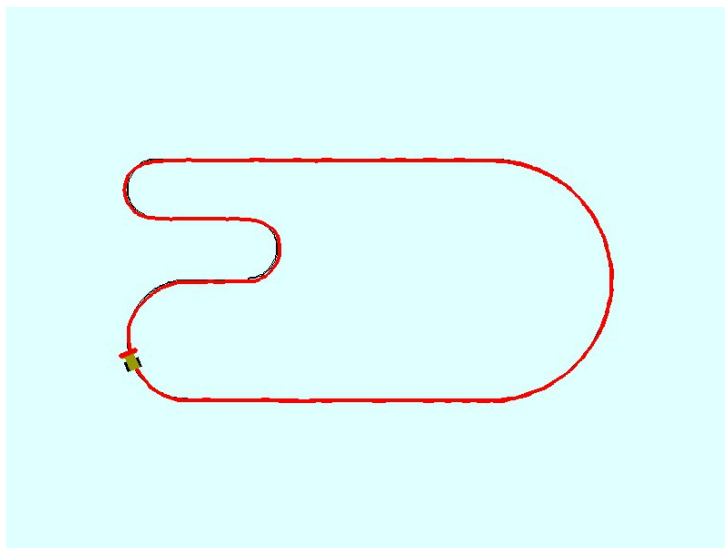


Figura 1: Trajetória do robô seguidor de linha após aprendizado com *Q-Learning*

O objetivo desses algoritmos neste contexto foi de avaliar métodos de aprendizado *on-policy* e *off-policy* no contexto do RL Livre de Modelo. Os algoritmos desenvolvidos resolvem o problema de Aprendizado por Reforço (*Reinforcement Learning - RL*) no caso em que o MDP é conhecido.

Os resultados foram então mostrados e discutidos neste relatório.

## 2 Introdução

### 2.1 Aprendizado *On* e *Off-Policy*

No contexto de algoritmos de aprendizado por reforço, um agente *on-policy* aprende o valor com base em sua ação atual  $\mathbf{a}$  derivada da política  $\pi$  atual, enquanto sua contraparte *off-policy* o aprende com base na ação  $\mathbf{a}^*$  obtida de uma outra política  $\mu$  [1]. No *Sarsa*, a avaliação de política é  $\epsilon$ -greedy. No *Q-learning*, essa política  $\pi$  é a política gananciosa  $\pi(s) = greedy(V_k(s))$  em relação ao estado  $s$  e à função valor  $V_k(s)$ [2].

### 2.2 Política $\epsilon$ -greedy

A ideia dessa abordagem é garantir exploração contínua no resultado, pois a ação tomada é a melhor até então (ação gulosa) com uma probabilidade de  $1 - \epsilon$ . Todas as demais ações terão igual probabilidade:

$$\pi(a|s) = \begin{cases} \frac{\epsilon}{m} + 1 - \epsilon, & a = \underset{a' \in A}{argmax} Q(s', a') \\ \frac{\epsilon}{m}, & \text{caso contrário} \end{cases} \quad (1)$$

### 2.3 Sarsa

A ideia do Sarsa é utilizar a recompensa  $R$  e o resultado do estado  $S'$  após se tomar a ação  $A'$  para recalculer o estado  $S$  dada a ação  $A$ :

$$Q(S, A) = Q(S, A) + \alpha(R + \gamma Q(S', A') - Q(S, A)) \quad (2)$$

Em que  $\alpha$  é a taxa de aprendizado e  $\gamma$  é o fator de desconto.

### 2.4 *Q-Learning*

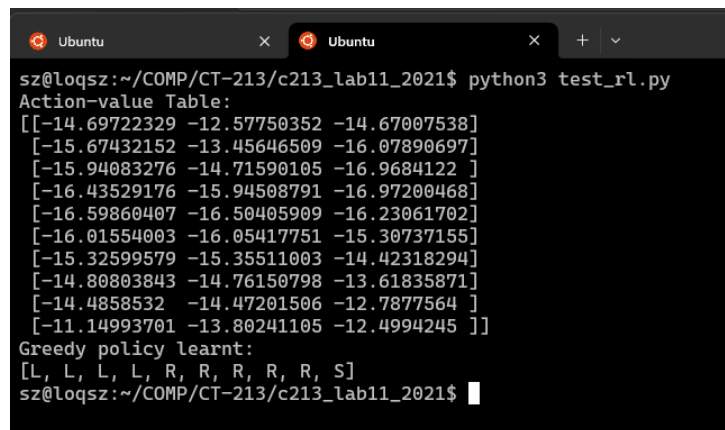
No *Q-Learning*, segue-se um aprendizado *off-policy*, ou seja, segue política de comportamento  $\mu$  enquanto aprende política alvo ótima  $\pi$ , em que  $\mu$  é  $\epsilon$ -greedy e  $\pi$  é greedy.

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \underset{a' \in A}{argmax} Q(S_{t+1}, a') - Q(S_t, A_t)) \quad (3)$$

## 3 Análise dos Resultados

### 3.1 Implementação dos algoritmos de RL

Após serem implementados os algoritmos Sarsa e *Q-Learning*, os algoritmos foram testados com o código *test\_rl.py*. O MDP de teste nesse *script* considera um “corredor” (“tabuleiro” unidimensional) de 10 células. As ações são STOP (ficar parado), LEFT (mover-se para esquerda) e RIGHT (mover-se para direita). As ações são consideradas determinísticas, no sentido de que sempre são executadas com total certeza. Além disso, esse MDP considera que ocorre “wrap” nos extremos do corredor: quando se executa LEFT na célula mais à esquerda, o agente surge na célula mais à direita. Analogamente, executar RIGHT na célula mais à direita faz o agente surgir na célula mais à esquerda. Finalmente, o estado objetivo é a célula mais à direita, de modo o agente recebe recompensa -1 em todas as células, exceto na célula objetivo, em que recebe recompensa 0. As figuras 2 e 3 mostram o resultado do teste desses algoritmos.



```
sz@loqsz:~/COMP/CT-213/c213_lab11_2021$ python3 test_rl.py
Action-value Table:
[[-14.69722329 -12.57750352 -14.67007538]
 [-15.67432152 -13.45646509 -16.07890697]
 [-15.94083276 -14.71590105 -16.9684122 ]
 [-16.43529176 -15.94508791 -16.97200468]
 [-16.59860407 -16.50405909 -16.23061702]
 [-16.01554003 -16.05417751 -15.30737155]
 [-15.32599579 -15.35511003 -14.42318294]
 [-14.80803843 -14.76150798 -13.61835871]
 [-14.4858532 -14.47201506 -12.7877564 ]
 [-11.14993701 -13.80241105 -12.4994245 ]]
Greedy policy learnt:
[L, L, L, L, R, R, R, R, R, S]
sz@loqsz:~/COMP/CT-213/c213_lab11_2021$
```

Figura 2: test\_rl para o algoritmo Sarsa

```
sz@loqsz:~/COMP/CT-213/c213_lab11_2021$ python3 test_rl.py
Action-value Table:
[[-1.14070342 -1.          -2.9701      ]
 [-2.00699329 -1.99       -3.940399   ]
 [-3.04436209 -2.9701     -4.90099464]
 [-3.99366142 -3.940399   -5.75956015]
 [-4.99686738 -4.8996782  -4.8996791  ]
 [-4.08391559 -4.07331223 -3.94039892]
 [-3.04048651 -3.00062591 -2.9701      ]
 [-2.09689383 -2.1643524  -1.99          ]
 [-1.05534158 -1.07906579 -1.          ]
 [ 0.          0.          -0.99         ]]
Greedy policy learnt:
[L, L, L, L, L, R, R, R, R, S]
sz@loqsz:~/COMP/CT-213/c213_lab11_2021$
```

Figura 3: test\_rl para o algoritmo *Q-Learning*

Aqui, vê-se que o *Q-Learning* apresenta um resultado mais intuitivo, ou seja, estados a  $X$  passos da casa alvo recebem uma recompensa de  $-X$ . Isso ocorre porque a atualização do estado considera uma política gulosa do estado seguinte, o que previne o algoritmo de fazer muita *exploration*. Em contrapartida, o *Sarsa* possui penalidades maiores para a função valor aprendida, já que a atualização do estado usa uma política  $\epsilon$ -greedy, que permite mais *exploration*.

### 3.2 Aprendizado da política do robô seguidor de linha

Após serem testados, os algoritmos foram utilizados para realizar o aprendizado do robô seguidor de linha. Ambos os algoritmos mostraram-se eficientes na tarefa a ser executada, gerando, com menos de 100 iterações, resultados muito próximos aos da figura 1. As figuras 4 e 5 mostram o gráfico de convergência para cada um dos algoritmos.

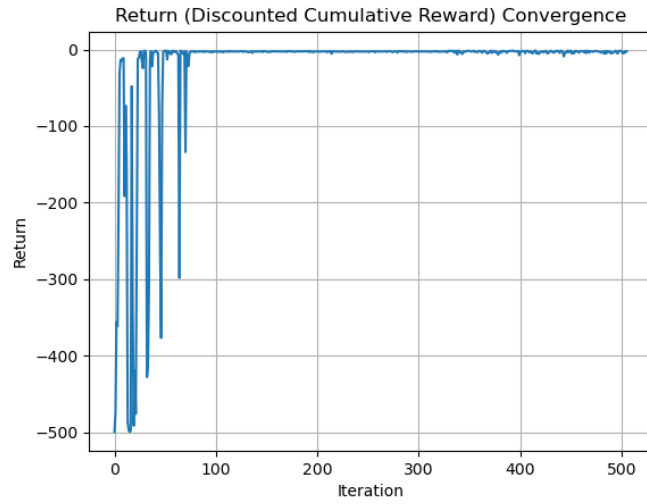


Figura 4: Convergência do algoritmo Sarsa para robô seguidor de linha

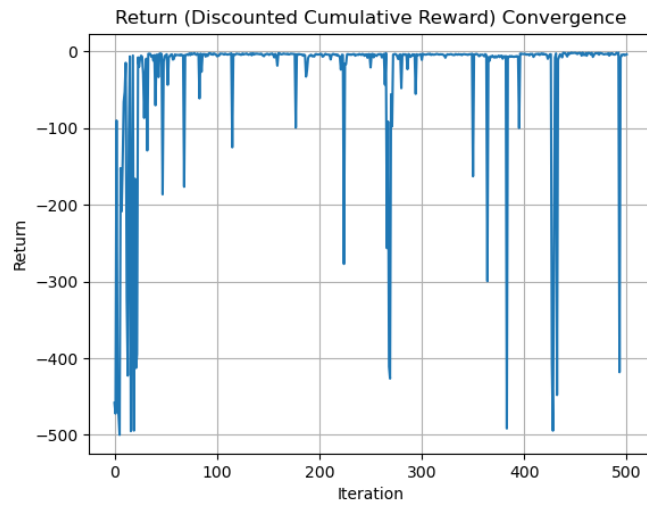


Figura 5: Convergência do algoritmo  $Q$ -Learning para robô seguidor de linha

Observa-se, pelo caráter de mais *exploration* do Sarsa, que, para este caso, ele apresenta resultados muito ruins (com recompensa na casa dos -400) nas 100 iterações iniciais, mas, no geral, após a convergência, tende a ter resultados mais consistentes, isto é, seguindo a linha mais fielmente. O  $Q$ -Learning, por sua vez, atinge resultados mais consistentes mais rapidamente (em 20 a 30 iterações), mas ocasionalmente gera resultados ruins.

## 4 Conclusão

O presente laboratório mostrou que algoritmos de *Reinforcement Learning* livres de modelo são capazes de ser utilizados na resolução de problemas quando não se tem o modelo nem os resultados da função ação-valor desse modelo. Métodos de aprendizado *on* e *off-policy* chegaram em uma solução próxima da ótima, porém, para este caso específico, o algoritmo Sarsa, de aprendizado *on-policy*, mostrou-se mais consistente em seus resultados.

Os algoritmos puderam ser treinados de forma a encontrar a trajetória ótima para o problema dado, com menos de 100 iterações, e chegando a resultados consistentes (recompensa praticamente igual a 0) na maioria das iterações após isso.

## Referências

- [1] Algoritmo de agente baseado em ia com reinforcement learning. Acessado em 19 de Maio de 2024. URL: <https://www.deeplearningbook.com.br/algoritmo-de-agente-baseado-em-ia-com-reinforcement-learning-parte-2/>.
- [2] Marcos Ricardo Omena de Albuquerque Maximo. Ct-213 - aula 12 - aprendizado por reforço livre de modelo. Apostila, 2020. Curso CT-213 - Inteligência Artificial para Robótica Móvel, Instituto Tecnológico de Aeronáutica.