



Universidade de Caxias do Sul
Nomes: Douglas Biazus
Professor(a): Carine Geltrudes Webber

Projeto Avaliativo
Aprendizado por Reforço:
Navegação em Labirintos com Q-Learning

Caxias do Sul
2025

Introdução

Como trabalho final de semestre, foi proposto implementar o algoritmo Q-Learning, com interface gráfica, para que um agente aprenda a navegar em um labirinto usando o algoritmo. Para o desenvolvimento do trabalho, foi utilizado a linguagem de programação Python, juntamente com a biblioteca Pygame para criação da visualização com interface gráfica. O programa executa múltiplas simulações para que o agente aprenda a encontrar o melhor caminho, recebendo pontuações positivas ao chegar no objetivo, mas perdendo pontuação ao passar por obstáculos, incentivando a encontrar o caminho mais eficiente.

Q-Learning

Dos algoritmos de aprendizado por reforço, há os que aprendem com o que é realizado pelo agente (*on-policy*) e os que aprendem como o que deveria ter sido feito (*off-policy*). O Q-learning é um algoritmo *off-policy*. Mesmo realizando uma ação não ótima, a tabela Q é atualizada com base na melhor ação para aquele estado, para que em próximas iterações escolhas melhores sejam feitas, buscando um melhor caminho no caso da navegação em labirintos.

0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				
0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				
0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				
0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				
0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				
0.0 0.0 -0.7 0.0	0.0 -0.9 -0.7 0.0	-140.0 0.0 0.0 0.0	0.0 0.0 0.0 0.0				

(Foto de parte da visualização em uma das primeiras iterações)

(Note que há pontuação negativa pelo agente ter se deslocado para cima, movimento não ótimo.)

0.0 0.0 0.0 0.0	0.0 -0.7 62.2 -0.7	0.0 43.3 79.0 -105.2	62.2 79.0 0.0 62.2	-120.0 100.0 26.2 0.0	100.0 126.2 59.0 0.0	126.2 159.0 0.0 0.0	0.0 0.0 0.0 0.0
0.0 0.0 -0.7 0.0	48.6 -1.4 -182.0 -1.1	62.2 149.0 0.7 37.6	79.0 62.2 0.0 48.7				
-0.7 0.0 27.2 -0.7	37.2 29.3 38.0 22.5	-150.3 38.0 48.7 29.4	62.2 48.7 0.0 38.0				
12.2 0.0 22.5 12.6	29.4 22.5 29.4 17.0	38.0 29.4 38.0 -176.5	48.7 38.0 0.0 29.4				
17.0 0.0 17.0 9.1	22.5 17.0 -176.5 12.6	29.4 176.0 29.4 17.0	38.0 29.4 0.0 22.5				
12.6 0.0 12.6 0.0	17.0 12.6 17.0 0.0	-176.5 17.0 22.5 0.0	29.4 22.5 0.0 0.0				

(Foto de parte da visualização na iteração final.) (2ª casa decimal ocultada)

(O agente prioriza a direção com maior valor na tabela Q.)

O algoritmo possui uma tabela em que os valores são atualizados conforme a fórmula:

$$Q[\text{estado}][\text{ação}] = Q[\text{estado}][\text{ação}] + \text{alfa} * (\text{recompensa} + \text{gama} * \max(Q[\text{próximo_estado}]) - Q[\text{estado}][\text{ação}])$$

- Alfa é a taxa de aprendizado, que determina o peso das novas iterações na Tabela Q. Um Alfa alto acelera o aprendizado do agente, fazendo ele dar mais peso a novas informações, porém podendo comprometer a qualidade da Tabela Q, enquanto que um Alfa baixo prioriza estabilidade, obrigando mais iterações e tornando o processo lento.
- Gama é o fator de desconto, com base no próximo estado (quadrado) que o agente estará presente. Um Gama alto faz o agente priorizar a recompensa futura. Um Gama baixo, faz priorizar a recompensa na posição atual.

Por fins de simplicidade, alfa e gama fixos foram utilizados no projeto. Testei o algoritmo diminuindo alfa em cada etapa, e usando alfa e gama mais altos, mas o resultado não sofreu alterações. O agente consegue alcançar o objetivo em menos de 3000 iterações independente dos valores testados, exceto ao diminuir o gama, porque um gama baixo reduz a importância da recompensa futura. No caso do labirinto atrapalhou na priorização do caminho mais curto.

Métricas

Em termos de métricas, por volta de 50 iterações, o agente já apresenta um bom desempenho, errando o caminho apenas pela aleatoriedade de 30% imposta. Essa aleatoriedade de 30% é necessária pois sem ela o agente poderia ficar preso em caminhos “vai-e-volta” nas primeiras iterações, não aprendendo. Foi necessário um mínimo de 18 iterações para que a Tabela Q mostrasse um caminho até o objetivo.

0.0	0.0	0.0	59.2	-200.8	-1.7	0.0	0.0
0.0 0.0	0.0 0.0	-0.9 -1.7	114.0 43.1	141.0 60.1	144.0 79.0	160.0 00.0	0.0 0.0
0.0	0.0	0.0	110.4	0.0	0.0	0.0	0.0
0.0	0.0	-0.9	127.8				
0.0 0.0	-0.9 -180.0	-180.0 -0.9	101.0 0.0				
0.0	0.0	0.0	-1.0				
0.0	-0.9	-180.0	114.0				
0.0 -0.9	-0.9 -0.9	-0.9 -0.9	-1.0 0.0				
0.0	0.0	0.0	-1.0				
-0.9	-0.9	-0.9	101.6				
0.0 -0.9	-1.0 -1.0	-0.9 -0.9	-1.0 0.0				
-0.9	-1.6	-180.7	-1.0				
-1.0	-1.7	-0.9	90.4				
0.0 -1.9	-1.7 -198.0	-198.0 1.0	61.0 0.0				
-1.0	-1.0	-0.9	53.3				
-1.0	-1.9	-180.0	80.3				
0.0 0.4	-1.0 1.5	-1.0 0.9	69.0 0.0				
0.0	0.0	0.0	0.0				

(Foto de parte da visualização na 18ª iteração.)

0.0	0.0	-0.9	114.0	-71.2	142.2	160.1	0.0
0.0 0.0	-0.9 -0.9	-0.9 -0.9	-1.0 43.1	143.0 60.1	160.0 79.0	179.0 00.0	0.0 0.0
0.0	0.0	0.0	114.0	0.0	0.0	0.0	0.0
0.0	-0.9	-0.9	127.8				
0.0 -0.9	-1.0 -180.0	-180.0 -0.9	114.0 0.0				
0.0	-1.6	0.0	101.6				
-0.9	-1.0	-180.0	114.0				
0.0 -1.7	-1.7 -1.7	-1.7 80.3	101.0 0.0				
-0.9	-0.9	-0.9	81.3				
-1.0	-1.0	-1.0	101.6				
0.0 -1.6	-1.8 -1.0	-1.0 81.3	90.0 0.0				
-1.0	-1.6	-132.9	80.4				
-1.0	-1.7	64.8	90.4				
0.0 -1.8	-1.7 -180.0	-180.0 0.4	80.0 0.0				
49.3	50.2	-0.9	71.3				
43.4	-1.9	-127.9	80.4				
0.0 0.4	55.0 3.2	63.0 1.4	71.0 0.0				
0.0	0.0	0.0	0.0				

(Foto de parte da visualização na 50ª iteração.)

Aplicações Práticas de Q-Learning

Por armazenar cada par de estado e ação em uma linha diferente, ou matriz diferente, o Q-Learning não é eficiente para aplicações com muitos dados. No jogo Five-in-a-Row, Gomoku, “cinco em linha” que tentei desenvolver, o número de estados possíveis e combinações no tabuleiro é gigante, obrigando a diminuir o número de casos. Só é possível de se fazer isso fazendo o Q-Learning de outras maneiras, como tratando linhas de “cuidado” e “ameaça de fim de jogo” ao invés de tratar o tabuleiro, limitando o Q-Learning e obrigando a criar mais funções com funcionamento do jogo para lidar com a recompensa. Não é possível criar uma punição e deixar o algoritmo aprender por conta o funcionamento do jogo igual é no labirinto. Se fossemos aplicar para um sistema de recomendação de filmes, também teríamos o mesmo problema: não dá para usar uma combinação de filmes assistidos e filmes novos para fazer recomendação com Q-Learning pelo volume de informações. Precisaríamos segmentar o problema em vários menores, fazendo um algoritmo para descobrir categorias de filmes (comédia, drama, ação, suspense) que o usuário tem interesse, outro para descobrir quais diretores de filmes recomendar filmes, e outro, por fim, para mostrar quais filmes as pessoas mais gostam em geral, e não individualmente. Simplesmente não é viável.

O Q-Learning é bom para resolver problemas simples, que dá pra dividir tudo em passos definidos, e em poucas ações. O jogo Snake em que a cobra precisa evitar bater nas paredes e buscar comida, e o jogo Pac-Man, que o agente precisa aprender a fugir dos fantasmas, ganhar pontos, pegar frutas, e terminar a fase o mais rápido, são dois exemplos de jogos que imagino serem boas aplicações do QLearning. Em controles de robôs, robótica? Talvez para evitar obstáculos, em mapas fixos (semelhante ao labirinto do projeto, um aspirador de pó evitando bater em objetos da casa) Em controle autônomo? Talvez para saber quais luzes ligar e desligar, levando em conta o horário e os hábitos dos moradores, para simular que há pessoas na casa mesmo quando a família está de viagem, aumentando a segurança. Fora desses cenários simples, o Q-Learning se torna muito limitado e não consegue lidar com ambientes complexos.