

# Inteligência Artificial

---

Prof. Fabrício Olivetti de França   Prof. Denis Fantinato

3º Quadrimestre de 2018

## **Aprendizado por Reforço**

Até esse instante assumimos que o agente tem um modelo completo do ambiente, sabe exatamente o estado atual, o conjunto de ações e as recompensas para cada estado.

No entanto, em muitos casos não temos conhecimento de nenhum desses.

Imagine jogar um jogo sem conhecer as regras e, após algumas centenas de jogadas, o seu oponente anuncia: *Você perdeu!*.

Vamos assumir um ambiente:

- **Totalmente observável:** o agente consegue conhecer o estado atual do ambiente.
- O agente não sabe em que cada ação resulta.
- O agente não sabe como o ambiente funciona.
- O efeito das ações são probabilísticas.

Isso torna impossível aplicar os algoritmos aprendidos na aula anterior para encontrar a política ótima (um processo de decisão de Markov desconhecido!).

Podemos tentar:

- Aprender uma **função utilidade** de cada estado e usar tal função para selecionar as ações que maximizam o valor total esperado.
- Aprender uma **função-Q**  $Q(s, a)$  que retorna um valor esperado de executar uma ação no estado atual.
- Aprender uma **função reflexiva** que retorna uma ação dado um estado.

# **Aprendizado Passivo**



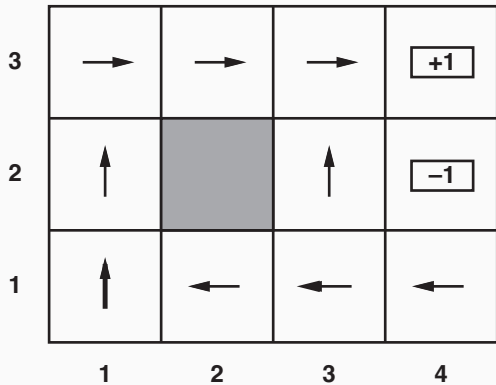
No **aprendizado passivo** desejamos estimar a função utilidade através de observações utilizando uma **política fixa**.

O objetivo é apenas determinar quão boa é uma política, estimando  $U^\pi$ .

O agente não conhece nem as transições  $P(s'|s, a)$  nem as recompensas  $R(s)$  para cada estado.

# Aprendizado Passivo

Considere a seguinte política:



O aprendizado passivo executa diversas **tentativas** de seguir a política até algum estado final.

## Aprendizado Passivo

Considere essas três tentativas:

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \\ \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \\ \mapsto (3, 2)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

$$(1, 1)_{-0.04} \mapsto (2, 1)_{-0.04} \mapsto (3, 1)_{-0.04} \mapsto (3, 2)_{-0.04} \mapsto (4, 2)_{-1}$$

Com isso temos uma amostra da utilidade de cada estado segundo a política  $\pi$ .

Para cada estado na amostra, pode-se estimar a utilidade como:

$$\sum_{t=0}^{\infty} \gamma^t R(S_t)$$

Observando a primeira tentativa:

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \\ \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

## Aprendizado Passivo

Assumindo  $R(S_t) = -0.04$  para estados não terminais e  $\gamma = 1$ , temos:

s	U
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(1,2)	0.84
(1,3)	0.88
(2,3)	0.92
(3,3)	0.96
(4,3)	1.00

Qual seria uma boa estimativa para  $U^\pi(1, 2)$ ?

s	U
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(1,2)	0.84
(1,3)	0.88
(2,3)	0.92
(3,3)	0.96
(4,3)	1.00



$$U^\pi(1, 2) = (0.76 + 0.84)/2 = 0.8$$

<u>s</u>	<u>U</u>
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(1,2)	0.84
(1,3)	0.88
(2,3)	0.92
(3,3)	0.96
(4,3)	1.00

Com a segunda tentativa:

$$\begin{aligned} (1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \\ \mapsto (3, 2)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1} \end{aligned}$$

Temos:

<u>s</u>	<u>U</u>
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(2,3)	0.84
(3,3)	0.88
(3,2)	0.92
(3,3)	0.96
(4,3)	1.00

Como atualizamos  $U^\pi(1, 2) = (0.76 + 0.84)/2 = 0.8$ ?

<u>s</u>	<u>U</u>
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(2,3)	0.84
(3,3)	0.88
(3,2)	0.92
(3,3)	0.96
(4,3)	1.00

Como atualizamos  $U^\pi(1, 2) = (0.76 + 0.84 + 0.76)/3 = 0.79$ ?

<u>s</u>	<u>U</u>
(1,1)	0.72
(1,2)	0.76
(1,3)	0.80
(2,3)	0.84
(3,3)	0.88
(3,2)	0.92
(3,3)	0.96
(4,3)	1.00

E na terceira tentativa:

$$(1, 1)_{-0.04} \mapsto (2, 1)_{-0.04} \mapsto (3, 1)_{-0.04} \mapsto (3, 2)_{-0.04} \mapsto (4, 2)_{-1}$$

Temos:

<u>s</u>	<u>U</u>
(1,1)	-1.16
(2,1)	-1.12
(3,1)	-1.08
(3,2)	-1.04
(4,2)	-1.00

O objetivo é obter:

$$U^\pi(s) = E_s \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$



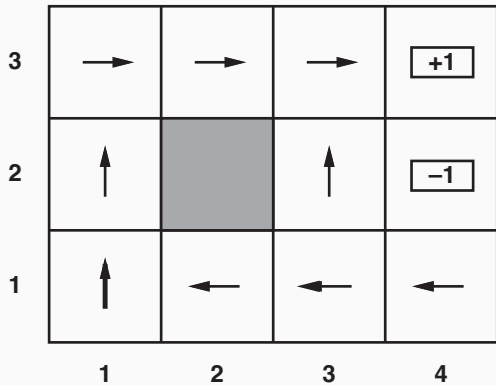
Uma forma simples para estimar a função de utilidade  $U^\pi(s)$  para um dado  $s$  é calcular a média da utilidade obtida em todas as ocorrências desse  $s$ .

Alternativamente, podemos utilizar a média móvel das  $n$  últimas ocorrências de  $s$ , dessa forma se nosso ambiente é alterado com o tempo, teremos sempre uma estimativa atual da utilidade.

Em um ambiente estacionário, com o número de amostras tendendo ao infinito, obtemos o valor real da utilidade para cada estado percorrido pela política  $\pi$ .

## Exercício

Dada a política:



## Exercício

E a amostra:

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \\ \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

Estime os valores da utilidade para cada estado utilizando a estimativa direta.

$$U(1, 1) = 0.72$$

$$U(1, 2) = 0.80$$

$$U(1, 3) = 0.84$$

$$U(2, 3) = 0.92$$

$$U(3, 3) = 0.96$$

$$U(4, 3) = 1.00$$

Note que essa forma de estimar a utilidade remete a um **aprendizado supervisionado** ou **indutivo**, em que temos uma amostra de valores de utilidade (a partir das recompensas) para cada estado e desejamos estimar um  $\hat{U}(s) \approx U^\pi(s)$ .

Na disciplina de **Aprendizado de Máquina** são apresentados diversos algoritmos desse tipo.



No entanto, a estimativa direta da utilidade não considera a noção de vizinhança dos estados (ela erra ao assumir que a utilidade dos estados é independente).

Uma outra forma de atualizar os valores de  $U^\pi$  é levar em conta as restrições da equação de Bellman.

$$U^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) U^\pi(s')$$

Porém,  $P(s'|s, \pi(s))$  não é conhecido: considera-se apenas uma transição e  $P(s'|s, \pi(s)) = 1$ .

## Aprendizado por Diferença Temporal

Partindo do algoritmo anterior, digamos que calculamos a estimativa de  $U^\pi(s)$  para os estados  $(1, 3)$  e  $(2, 3)$  após a primeira amostra:

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \\ \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

$$U^\pi(1, 3) = 0.84$$

$$U^\pi(2, 3) = 0.92$$

Pela equação de Bellman, temos que a utilidade do estado  $(1, 3)$  deveria obedecer ( $\gamma = 1$ ):

$$U^\pi(1, 3) = -0.04 + U^\pi(2, 3)$$

E no nosso caso:

$$U^\pi(1, 3) = -0.04 + 0.92 = 0.88$$

Sendo que a estimativa direta é 0.84.

Podemos então atualizar  $U^\pi$  como:

$$U^\pi(s) = U^\pi(s) + \alpha(R(s) + \gamma U^\pi(s') - U^\pi(s))$$

com  $\alpha$  sendo a taxa de aprendizado e  $0 \leq \alpha \leq 1$ .

Em nosso exemplo, seria:

$$U^\pi(1, 3) = U^\pi(1, 3) + \alpha(-0.04 + 0.92 - U^\pi(1, 3))$$

Que para um  $\alpha = 1$ :

$$U^\pi(1, 3) = 0.84 + (-0.04 + 0.92 - 0.84) = 0.88$$

Mas para  $\alpha = 0.1$ :

$$U^\pi(1, 3) = 0.84 + 0.1(-0.04 + 0.92 - 0.84) = 0.844$$

Ele atualiza um pouquinho na direção apontada pela equação de Bellman.



Note que a atualização do estado  $s$  envolve apenas o estado seguinte  $s'$  que foi amostrado.

Lembrando que dada uma ação  $a$  podemos chegar em diferentes estados  $s'$  probabilisticamente, um estado  $s'$  raro (e que afeta negativamente o sistema) poderia subestimar ou superestimar o valor da utilidade.

Mas, como trabalhamos com amostragens e essa equação leva ao cálculo de um valor médio, um evento raro afeterá pouco no resultado.

Se temos um  $\alpha$  inversamente proporcional a quantas vezes um estado  $s$  foi visitado (quanto maior o número de visitas, menor o  $\alpha$ ), garantimos a convergência para o valor correto de  $U^\pi(s)$ .

Um exemplo é fazer:

$$\alpha(n) = \frac{60}{59+n}$$

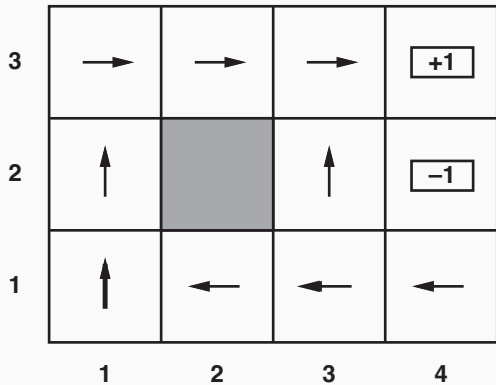
sendo  $n$  o número de vezes que o estado foi visitado.

Se temos um espaço de busca com muitos estados, o algoritmo pode se tornar intratável.

Uma alternativa é utilizar uma heurística para atualizar apenas determinados estados.

# Exercício

Dada a política:



## Exercício

E as amostras:

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \\ \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

$$(1, 1)_{-0.04} \mapsto (1, 2)_{-0.04} \mapsto (1, 3)_{-0.04} \mapsto (2, 3)_{-0.04} \mapsto (3, 3)_{-0.04} \\ \mapsto (3, 2)_{-0.04} \mapsto (3, 3)_{-0.04} \mapsto (4, 3)_{+1}$$

## Exercício

Tome como estimativa inicial os valores obtidos pela estimativa direta do exercício anterior com a primeira amostra.

Usando a segunda amostra, estime os valores da utilidade pela diferença temporal (com  $\alpha = 0.1, \gamma = 1$ ).



$$U(1, 1) = 0.72$$

$$U(1, 2) = 0.80$$

$$U(1, 3) = 0.84$$

$$U(2, 3) = 0.92$$

$$U(3, 3) = 0.96$$

$$U(4, 3) = 1.00$$

$$U(1, 1) = 0.72 + 0.1(-0.04 + 0.8 - 0.72) = 0.724$$

$$U(1, 2) = 0.80 + 0.1(-0.04 + 0.84 - 0.8) = 0.80$$

$$U(1, 3) = 0.844$$

$$U(2, 3) = 0.92$$

$$U(3, 3) = 0.96$$

$$U(4, 3) = 1.00$$

## **Aprendizado por Reforço Ativo**

Nos algoritmos anteriores focamos em estimar a função de utilidade para uma política fixa, assumindo que ela era ótima.

No Aprendizado Ativo, queremos também estimar a política ótima.

Basicamente temos:

```
def ativo(U, pi):  
    U = estima(U, pi)  
    pi = politica(U)  
    return U, pi
```

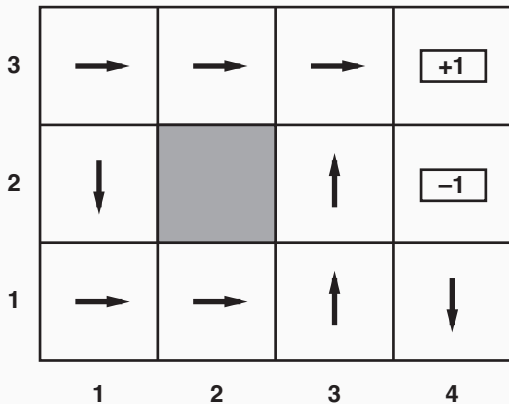
Ou seja, dados  $U$ ,  $\pi$ , fazemos uma nova simulação para aprimorar a estimativa de  $U$  (através dos algoritmos anteriores) e, então, estimamos um novo  $\pi$ , baseado no  $U$  atual (iteração de política).

O uso desse algoritmo leva a criação de um **agente guloso** pois ele sempre irá caminhar pelo melhor caminho encontrado até o momento.



## Estimando a utilidade e a política

Considere que em certo momento o agente encontrou a seguinte política que o levou até a posição  $+1$ :



Como essa foi a melhor solução encontrada até então, o agente passará a **reforçar** esse caminho com poucas alterações causadas pela probabilidade de execução de ação.

Difícilmente ele irá encontrar a nossa solução ótima uma vez que ele entra nessa auto-alimentação da solução atual.

Portanto, o agente guloso tem um problema em que ele não explora novos caminhos caso tenha encontrado algo *bom*.

Embora a estimativa de  $U$  seja ótima segundo a política  $\pi$  atual, não necessariamente a política é ótima.

## Exploração vs Exploração

Uma forma de evitar isso é balancear a **exploração** e a **exploração**.

## Exploração vs Exploração

- **Exploração:** percorrer caminhos pouco (ou ainda não) percorridos.
- **Exploração:** reforçar o melhor caminho encontrado até então.

## Exploração vs Exploração

Um bom equilíbrio entre exploração e exploração pode determinar o sucesso de uma heurística.

Uma ideia simplista é favorecer caminhos não explorados até que estes tenham sido explorados por um  $n$  número de vezes.



Alguns algoritmos alternam entre exploração e exploração probabilisticamente.

# **Aprendizado Ativo com Diferença Temporal**

Q-Learning é um algoritmo de diferença temporal que aprende uma função  $Q(s, a)$  que indica a qualidade em escolher a ação  $a$  no estado  $s$ .

Essa função é relacionada com a função utilidade de tal forma que:

$$U(s) = \max_a Q(s, a)$$

Similar a diferença temporal, no algoritmo Q-Learning atualizamos a função  $Q(s, a)$  como:

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

que é executado toda vez que aplicamos uma ação  $a$  em  $s$  levando ao estado  $s'$ .

Como não precisamos estimar  $P(s' | s, a)$  e dizemos que Q-Learning é um algoritmo **livre de modelo** (*model free*).

```
def qlearn(percept, state, goals, A, gamma, alpha):
    sn, rn          = percept
    s, a, r, N, Q = state

    # se for estado final, usa a recompensa
    if sn in goals:
        Q[(sn, None)] = rn
```

```
# atualiza Q  
else:  
    N[(s,a)] += 1  
    maxQ      = max(Q[(sn,ai)] for ai in A[sn])  
    Q[(s,a)] = Q[(s,a)] + alpha * (r +  
                                     gamma*maxQ - Q[(s,a)])
```



```
# avalia o proximo estado e  
# a acao que deve ser tomada  
qvals    = [(Q[(sn,ai)], ai) for ai in A[sn]]  
an       = argmax(qvals)  
state    = (sn, an, rn, N, Q)  
  
return an, state
```

Sendo o estado  $s$  atual um estado final, atualiza o valor de  $Q(s, \cdot)$  para a recompensa final. Note que não temos uma ação a ser executada.

Caso não seja um estado final, atualiza o valor de  $Q(s, a)$  para o novo valor utilizando a diferença temporal.

Em seguida, escolhe a próxima ação que maximize  $Q(s', a)$ .

Um problema do Q-Learning é que ele não leva em conta as consequências da ação atual, ou seja, ele assume que a ação escolhida para executar terá consequências futuras ótimas.

## SARSA: State-Action-Reward-State-Action

Uma forma de aliviar tal problema é modificando a equação de atualização de  $Q(s, a)$  para:

$$Q(s, a) = Q(s, a) + \alpha(R(s) + \gamma Q(s', a') - Q(s, a)),$$

ou seja, além do estado e ação atual  $s, a$ , também já deve ser conhecidos o estado e ação futura  $s', a'$ . Basta alterar a linha pertinente no algoritmo anterior.

Q-Learning não usa política, pois escolhe apenas o melhor valor Q (ou seja, é *off-policy*).

SARSA usa política, pois utiliza estado e ação futura escolhidos pela política (ou seja, é *on-policy*).

## Exploração vs Exploração

- **Exploração:** aprender possíveis valores de  $Q(s, a)$  para combinações de estado e ação ainda não observados.
- **Exploração:** seguir um caminho mais provável de maximizar a recompensa, dado o conhecimento atual.



## Exploração vs Exploração

Com certa probabilidade escolhe entre explorar e explorar. Se escolher explorar, o agente segue a ação que maximiza o valor de  $Q(s', a)$ .

Caso seja escolhido explorar, ou escolhe-se uma ação completamente aleatória, ou uma dentre aquelas que foram pouco exploradas.

Além disso, podemos fazer um passo de aprendizado adaptativo com  $\alpha = \alpha_0/t$  sendo  $t$  o número de episódios jogados até então.

De todo modo, o algoritmo Q-Learning ainda sofre em ambientes contendo muitos possíveis estados. Esse algoritmo costuma ser mais lento do que aprendizado passivo, porém permite adaptar facilmente para os casos em que não temos um ambiente como o requisitado pelo MDP.

Uma alternativa para tratar o problema de espaço de busca muito grande, é tentar encontrar uma função fechada para  $Q(s, a)$  que não necessite de uma tabela com as informações dos valores  $Q$ .

## **Generalizando para Estados não vistos**

Uma outra forma de pensarmos na função utilidade é de ela depender de uma variável  $\theta$  de dimensão menor que a representação dos estados.

Dessa forma desejamos obter  $\hat{U}_\theta(s) \approx U(s)$ .

Seguindo nosso exemplo, poderíamos pensar em usar as coordenadas  $[x, y]$  do mundo de nosso agente e obter:

$$\hat{U}_\theta(x, y) = \theta_0 f_0(s) + \theta_1 f_1(s) + \dots + \theta_n f_n(s)$$

$$\hat{U}_\theta(x, y) = \theta_0 + \theta_1 x + \theta_2 y$$

Assim, se  $(\theta_0, \theta_1, \theta_2) = (0.5, 0.2, 0, 1)$ , então  $\hat{U}_\theta(1, 1) = 0.8$ .



Dada uma amostra de tentativas, codificamos cada estado  $s$  em características  $([x, y])$  e, então, utilizamos um algoritmo de aprendizado supervisionado para calcular os valores de  $\theta$  na função  $\hat{U}_\theta(s)$ .

## Estimando uma função

```
def utilityModel(nTrials, S,R,s, goals, gamma, nextState):  
    dataX, dataY = makeNTrials(nTrials)  
  
    lr = LinearRegression()  
    lr.fit(dataX, dataY)  
  
    model = lambda s: lr.predict([[1, s[0], s[1]]])[0]  
  
    return model
```

Para o nosso problema com 1000 execuções, obtemos a seguinte equação:

$$\hat{U}_\theta(x, y) = \theta_0 + \theta_1 x + \theta_2 y = -0.06 + 0.11x + 0.17y$$

Que se traduz em:

0.56	0.68	0.79	0.90
0.40		0.62	0.74
0.23	0.34	0.46	0.57

E na política:

->	->	->	0.90
^		^	0.74
^	->	^	^

## Estimando uma função

Seja agora uma nova característica  $f_3(s) = (x - x_g)^2 + (y - y_g)^2$ , em que  $x_g$  e  $y_g$  são as coordenadas do objetivo. A utilidade fica:

0.54	0.51	0.50	0.48
0.38		0.34	0.32
0.23	0.20	0.18	0.17

E na política:

->	<-	->	0.48
^		^	0.32
^	<-	^	^

Essa forma de encontrar uma política ótima tem a vantagem de codificar o espaço de estados, que pode ser bem numeroso, para um conjunto de estados gerenciáveis.



Além disso, por estimar uma função, os modelos de aprendizado conseguem generalizar para estados que ainda não foram observados através de interpolação.

## **Busca por uma Política**

Uma última alternativa para estimar a política ótima é fazer a busca no **espaço de políticas**.

A ideia geral é seguir um algoritmo de busca em que, dada uma política inicial, altera essa política até que nenhuma outra melhora possa ser feita.

Um possível algoritmo é o **Hill Climbing** que parte de uma solução inicial  $s_0$ , testa todas as soluções vizinhas e substitui a solução atual pelo melhor vizinho, caso tenha um melhor desempenho.

A busca para no momento que não existir um vizinho melhor que a solução inicial.

```
def hillClimbing(mdp, s0, goals, nTrials):  
    S, A, R, P, gamma = mdp  
  
    piCurrent = { s : np.random.choice(A[s]) for s in S }  
    Ucurrent  = evaluate(mdp, piCurrent, nTrials)  
    improved  = True
```

```
while improved:
    improved = False
    bestNeigh, Uneigh = genNeighbors(mdp, piCurrent, s0,
                                     goals, nTrials)

    if Uneigh[s0] > Ucurrent[s0]:
        improved = True
        Ucurrent = deepcopy(Uneigh)
        piCurrent = deepcopy(bestNeigh)
return piCurrent
```

Os algoritmos de Aprendizado por Reforço devem ser aplicados quando não temos todas as informações do problema de busca e queremos uma estimativa da função de utilidade dos estados ou da política ótima.



Esses algoritmos podem ser de:

- **Aprendizado Passivo:** em que fixamos a informação de política para estimar a utilidade.
- **Aprendizado Ativo:** alternamos entre estimar a política e a função utilidade.
- **Aprendizado de Modelo:** quando utilizamos um modelo de regressão ou classificação em uma amostra de jogadas.
- **Aprendizado por Busca:** em que buscamos por uma política ótima utilizando uma estimativa da utilidade como função-objetivo.