

Relatório do Laboratório 10 - Programação Dinâmica

1. Breve Explicação em Alto Nível da Implementação

A implementação foi realizada baseada na execução de 3 funções principais: *policy_evaluation*, *value_iteration* e *policy_iteration*.

1.1. Avaliação de Política

A função *policy_evaluation* tinha por objetivo avaliar a função valor v_π para determinada política π . Para isso, deve-se resolver a equação de Bellman, encontrando a tabela que mapeia a política a cada um dos estados e a cada uma das ações.

Utilizando programação dinâmica, pode-se resolver a equação iterativamente. A partir de um chute inicial, um valor v_{k+1} é encontrado para a iteração $k + 1$ a partir do valor de v_k . Com $k \rightarrow \infty$, a função tende a convergir para $v_k = v_\pi$.

Definindo uma condição de parada de número de iterações máximo ou de uma diferença máxima entre funções valor consecutivas ϵ , podemos definir um suficientemente próximo para v_π . O pseudocódigo implementado, retirado do livro do Sutton, é demonstrado abaixo.

```
1. Initialization
    $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$ ;  $V(\text{terminal}) \doteq 0$ 

2. Policy Evaluation
   Loop:
      $\Delta \leftarrow 0$ 
     Loop for each  $s \in \mathcal{S}$ :
        $v \leftarrow V(s)$ 
        $V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$ 
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
   until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)
```

1.2. Iteração de Valor

A iteração de valor, diferentemente da iteração de política, permite que a melhoria da política ocorra juntamente com a avaliação da mesma, em uma passagem só. O algoritmo funciona por tornar a equação a função de otimalidade de Bellman em uma regra de atualização, pegando o valor máximo do estado para cada uma das ações. O pseudocódigo implementado, retirado do livro do Sutton, é demonstrado abaixo.

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

| $\Delta \leftarrow 0$

| Loop for each $s \in \mathcal{S}$:

| $v \leftarrow V(s)$

| $V(s) \leftarrow \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

| $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$\pi(s) = \arg \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$

1.3. Iteração de Política

A iteração de política é um algoritmo para a identificação da política ótima a partir de uma sequência de melhorias monotônicas. A partir de uma política π pode-se encontrar sua função valor v_π , e essa função pode ser então utilizada para encontrar uma política melhor

π' . Isso é garantido pois como uma MDP tem um número finito de políticas determinísticas, esse processo deve convergir para uma política ótima em um número finito de iterações.

O algoritmo funciona em duas etapas: avaliação e melhoria. Na etapa de avaliação, a função valor v_k é identificada a partir de uma política π com método de avaliação de política definido em 1.1. A partir dessa função, a próxima melhor política π' é retirada a partir de uma política gulosa. Essa política π é então usada para encontrar outra função de valor v_{k+1} , repetindo o ciclo. O ciclo é interrompido com uma condição de iterações máximas ou até que $\max |v_k - v_{k+1}| < \epsilon$. O pseudocódigo implementado, retirado do livro do Sutton, é demonstrado abaixo.

1. Initialization
 $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$; $V(\text{terminal}) \doteq 0$
2. Policy Evaluation
 Loop:
 $\Delta \leftarrow 0$
 Loop for each $s \in \mathcal{S}$:
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)
3. Policy Improvement
 $\text{policy-stable} \leftarrow \text{true}$
 For each $s \in \mathcal{S}$:
 $\text{old-action} \leftarrow \pi(s)$
 $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$
 If $\text{old-action} \neq \pi(s)$, then $\text{policy-stable} \leftarrow \text{false}$
 If policy-stable , then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

2. Tabelas Comprovando Funcionamento do Código

2.1. Caso $p_c = 1,0$ e $\gamma = 1,0$

2.1.1. Avaliação de Política

```
Evaluating random policy, except for the goal state, where policy always executes stop:
Value function:
[ -384.09, -382.73, -381.19, * , -339.93, -339.93]
[ -380.45, -377.91, -374.65, * , -334.92, -334.93]
[ -374.34, -368.82, -359.85, -344.88, -324.92, -324.93]
[ -368.76, -358.18, -346.03, * , -289.95, -309.94]
[ * , -344.12, -315.05, -250.02, -229.99, * ]
[ -359.12, -354.12, * , -200.01, -145.00, 0.00]
Policy:
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , SURDL , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ * , SURDL , SURDL , SURDL , SURDL , * ]
[ SURDL , SURDL , * , SURDL , SURDL , S ]
```

2.1.2. Iteração de Valor

```

Value iteration:
Value function:
[ -10.00, -9.00, -8.00, * , -6.00, -7.00]
[ -9.00, -8.00, -7.00, * , -5.00, -6.00]
[ -8.00, -7.00, -6.00, -5.00, -4.00, -5.00]
[ -7.00, -6.00, -5.00, * , -3.00, -4.00]
[ * , -5.00, -4.00, -3.00, -2.00, * ]
[ -7.00, -6.00, * , -2.00, -1.00, 0.00]
Policy:
[ RD , RD , D , * , D , DL ]
[ RD , RD , D , * , D , DL ]
[ RD , RD , RD , R , D , DL ]
[ R , RD , D , * , D , L ]
[ * , R , R , RD , D , * ]
[ R , U , * , R , R , SURD ]

```

2.1.3. Iteração de Política

```

Policy iteration:
Value function:
[ -10.00, -9.00, -8.00, * , -6.00, -7.00]
[ -9.00, -8.00, -7.00, * , -5.00, -6.00]
[ -8.00, -7.00, -6.00, -5.00, -4.00, -5.00]
[ -7.00, -6.00, -5.00, * , -3.00, -4.00]
[ * , -5.00, -4.00, -3.00, -2.00, * ]
[ -7.00, -6.00, * , -2.00, -1.00, 0.00]
Policy:
[ RD , RD , D , * , D , DL ]
[ RD , RD , D , * , D , DL ]
[ RD , RD , RD , R , D , DL ]
[ R , RD , D , * , D , L ]
[ * , R , R , RD , D , * ]
[ R , U , * , R , R , SURD ]

```

2.2. Caso $p_c = 0,8$ e $\gamma = 0,98$

2.2.1. Avaliação de Política

```

Evaluating random policy, except for the goal state, where policy always executes stop:
Value function:
[ -47.19, -47.11, -47.01, * , -45.13, -45.15]
[ -46.97, -46.81, -46.60, * , -44.58, -44.65]
[ -46.58, -46.21, -45.62, -44.79, -43.40, -43.63]
[ -46.20, -45.41, -44.42, * , -39.87, -42.17]
[ * , -44.31, -41.64, -35.28, -32.96, * ]
[ -45.73, -45.28, * , -29.68, -21.88, 0.00]
Policy:
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , SURDL , SURDL , SURDL ]
[ SURDL , SURDL , SURDL , * , SURDL , SURDL ]
[ * , SURDL , SURDL , SURDL , SURDL , * ]
[ SURDL , SURDL , * , SURDL , SURDL , S ]

```

2.2.2. Iteração de Valor

```

Value iteration:
Value function:
[ -11.65, -10.78, -9.86, * , -7.79, -8.53]
[ -10.72, -9.78, -8.78, * , -6.67, -7.52]
[ -9.72, -8.70, -7.59, -6.61, -5.44, -6.42]
[ -8.70, -7.58, -6.43, * , -4.09, -5.30]
[ * , -6.43, -5.17, -3.87, -2.76, * ]
[ -8.63, -7.58, * , -2.69, -1.40, 0.00]
Policy:
[ D , D , D , * , D , D ]
[ D , D , D , * , D , D ]
[ RD , D , D , R , D , D ]
[ R , RD , D , * , D , L ]
[ * , R , R , D , D , * ]
[ R , U , * , R , R , S ]
-----

```

2.2.3. Iteração de Política

```

Policy iteration:
Value function:
[ -11.65, -10.78, -9.86, * , -7.79, -8.53]
[ -10.72, -9.78, -8.78, * , -6.67, -7.52]
[ -9.72, -8.70, -7.59, -6.61, -5.44, -6.42]
[ -8.70, -7.58, -6.43, * , -4.09, -5.30]
[ * , -6.43, -5.17, -3.87, -2.76, * ]
[ -8.63, -7.58, * , -2.69, -1.40, 0.00]
Policy:
[ D , D , D , * , D , D ]
[ D , D , D , * , D , D ]
[ R , D , D , R , D , D ]
[ R , D , D , * , D , L ]
[ * , R , R , D , D , * ]
[ R , U , * , R , R , S ]
-----

```

3. Discussão dos Resultados

A partir dos resultados encontrados podemos perceber características interessantes sobre cada um dos testes realizados. Para o teste em que $p_c = 1$ e $\gamma = 1$ percebemos que a função valor foi maximizada em comparação com a função valor com política aleatória. Percebe-se que essa política limita as ações possíveis, encontrando caminhos ótimos a serem percorridos. Percebe-se também que o método de iteração de valor obteve a mesma função de valor final ao método de iteração de política, resultado já esperado.

Já para o teste em que $p_c = 0,8$ e $\gamma = 0,98$, percebe-se uma restrição maior da política com relação às ações possíveis para cada um dos estados. Isso se dá pelo fato de que os parâmetros p_c e γ induzem a iteração da função valor à resultados que favorecem recompensas mais imediatas, fazendo com que menos caminhos possíveis possam ser percorridos.