

**Relatório do Laboratório 10 - Programação Dinâmica**

**1. Breve Explicação em Alto Nível da Implementação**

**1.1. Avaliação de Política**

Avaliação política é uma função que, dada uma política, calcula os valores de todos os estados possíveis. Ela foi implementada utilizando a equação de bellman de expectativa, segundo a qual:

$$v_{\pi}(s) = \sum_{a \in A} \pi(a|s) r(s, a) + \gamma \sum_{a \in A} \sum_{s' \in S} \pi(a|s) p(s'|s, a) v_{\pi}(s')$$

Assim, utilizando iteração de valor síncrona, ou seja, só atualizando a função valor após o fim da iteração, o valor era atualizado até uma relativa convergência, em que a maior diferença entre o iterado anteriormente e o novo deve ser menor que um limite epsilon.

**1.2. Iteração de Valor**

Na iteração de valor, para cada estado, se seleciona o valor segundo a ação que retorna o maior valor. Em outras palavras, o valor de cada estado pode ser representado pela seguinte fórmula:

$$v_{k+1}(s) = \max_{a \in A} \left( r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) v_k(s') \right)$$

Assim, teremos uma tabela de valores formada e a partir dela pode-se determinar a política ótima por uma política gulosa("greedy").

**1.3. Iteração de Política**

Na iteração política, inicialmente começamos com uma política e uma tabela de valores aleatórias. Depois disso, entra-se num loop com um número de iterações suficientemente grandes para garantir convergência em que sucessivamente se atualiza o valor segundo a avaliação política e depois, a partir da tabela de valores, se estabelece uma nova política greedy, que será utilizada para futuramente atualizar o valores.

**2. Tabelas Comprovando Funcionamento do Código**

Basta colocar os *prints* das tabelas (ou cópias das saídas do programa)

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

**2.1.1. Avaliação de Política**

```

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

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 ]
[ RD , D , D , R , D , D ]
[ R , RD , D , * , D , L ]
[ * , R , R , D , D , * ]
[ R , U , * , R , R , S ]
-----

```

## 3. Discussão dos Resultados

Vemos que o resultado a value function dados  $\rho$  e  $\gamma$  devem ser os mesmos para ambos os métodos, value iteration e policy iteration. Vê-se que a policy evaluation, mesmo usando uma política random, leva ao resultado final, tendo o aumento a função policy para chegar lá.