

## Processo de Decisão de Markov

1

## Percepção e Ação

O objetivo final de um algoritmo de controle robótico é escolher as ações corretas para executar uma tarefa determinada

Problemas:

- ❑ Incerteza na percepção
- ❑ Incerteza na ação

SCE-5868 Denis F. Wolf

2

## Incerteza nas ações de um robô

Situações possíveis:

Ações determinísticas  
X  
Ações estocásticas

Ambiente parcialmente observável  
X  
Ambiente totalmente observável

SCE-5868 Denis F. Wolf

3

## Incerteza nas ações de um robô

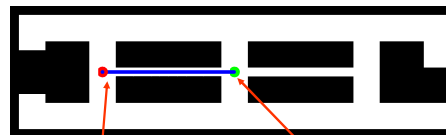
- ❑ Diversos resultados teóricos na robótica consideram determinístico o efeito das ações executadas
- ❑ Na maioria das situações reais, é insuficiente planejar uma seqüência de ações e executá-las "cegamente".
- ❑ É necessário considerar a incerteza quando se planeja uma seqüência de ações.
- ❑ Estratégia: antecipar a incerteza e planejar ações para qualquer situação possível

SCE-5868 Denis F. Wolf

4

## Situação ideal

Ações determinísticas e ambiente totalmente observável

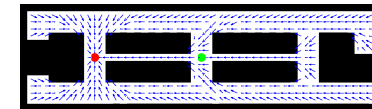


Posição de destino do robô      Posição inicial do robô

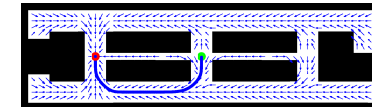
SCE-5868 Denis F. Wolf

## Plano global de ações

Ações determinísticas e ambiente totalmente observável



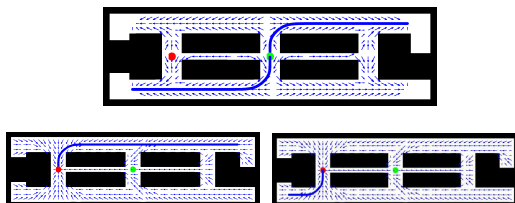
Ações estocásticas e ambiente totalmente observável



SCE-5868 Denis F. Wolf

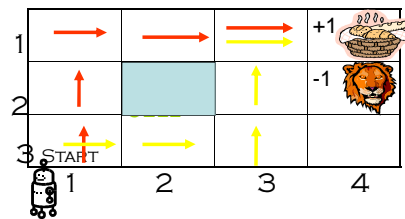
## Plano global de ações

Ações estocásticas e ambiente parcialmente observável



SCE-5868 Denis F. Wolf

Qual o melhor plano para chegar a recompensa?



SCE-5868 Denis F. Wolf

## Processo de decisão de Markov PDM

□ **Dados:**

- Estados  $x$
- Ações  $u$
- Probabilidades de transição  $p(x'|u,x)$
- Função de Recompensa  $r(x,u)$

□ **Obter:**

- Plano  $\pi(x)$  que maximiza a recompensa esperada

D1

SCE-5868 Denis F. Wolf

## Recompensas e planos

- Plano (estado observável):  $\pi: x_t \rightarrow u_t$
- Plano (caso geral):  $\pi: z_{t-1}, u_{t-1} \rightarrow u_t$
- Recompensa esperada:  $R_T = E \left[ \sum_{\tau=1}^T \gamma^\tau r_{t+\tau} \right]$ 
  - $T=1$ : plano greedy
  - $T>1$ : horizonte finito,  $\gamma=1$
  - $T=\text{infinity}$ : horizonte infinito, tipicamente  $\gamma < 1$

SCE-5868 Denis F. Wolf

## Recompensas e planos

- Plano ótimo para 1 passo: 
$$\pi_1(x) = \underset{u}{\operatorname{argmax}} r(x,u)$$
- Função de valor para o melhor plano de 1 passo: 
$$V_1(x) = \gamma \max_u r(x,u)$$

SCE-5868 Denis F. Wolf

## Plano de 2 etapas

- Plano ótimo: 
$$\pi_2(x) = \underset{u}{\operatorname{argmax}} \left[ r(x,u) + \int V_1(x') p(x'|u,x) dx' \right]$$
- Função de valor: 
$$V_2(x) = \gamma \max_u \left[ r(x,u) + \int V_1(x') p(x'|u,x) dx' \right]$$

SCE-5868 Denis F. Wolf

**Slide 9**

---

**D1**

Denis; 03/06/2009

## Plano de T etapas

- Plano ótimo:

$$\pi_T(x) = \operatorname{argmax}_u \left[ r(x,u) + \int V_{T-1}(x') p(x'|u,x) dx' \right]$$

- Função de valor:

$$V_T(x) = \gamma \max_u \left[ r(x,u) + \int V_{T-1}(x') p(x'|u,x) dx' \right]$$

SCE-5868 Denis F. Wolf

## Infinitas etapas

- Plano ótimo (equação de Bellman):

$$V_\infty(x) = \gamma \max_u \left[ r(x,u) + \int V_\infty(x') p(x'|u,x) dx' \right]$$

SCE-5868 Denis F. Wolf

## Value Iteration

- for all  $x$  do

$$\hat{V}(x) \leftarrow r_{\min}$$

- endfor

- repeat until convergence

- for all  $x$  do

$$\hat{V}(x) \leftarrow \gamma \max_u \left[ r(x,u) + \int \hat{V}(x') p(x'|u,x) dx' \right]$$

- endfor

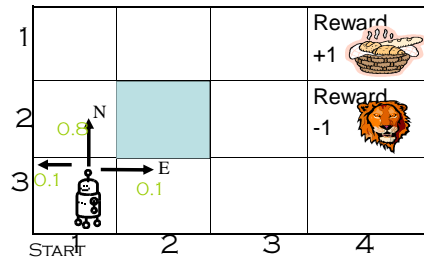
- endrepeat

$$\pi(x) = \operatorname{argmax}_u \left[ r(x,u) + \int \hat{V}(x') p(x'|u,x) dx' \right]$$

SCE-5868 Denis F. Wolf

## Exemplo: Qual o melhor plano de ação?

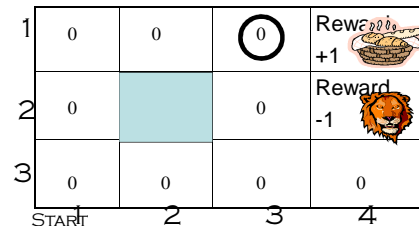
- O robô pode andar nas seguintes direções: N, E, S, W
- Cada ação tem um custo de  $-1/25$



SCE-5868 Denis F. Wolf

## 1ª Iteração: célula (3,1)

- East:  $-1/25 + [0.8 * 1 + 0.1 * 0 + 0.1 * 0] = 0.76$
- North/South:  $-1/25 + [0.8 * 0 + 0.1 * 1 + 0.1 * 0] = 0.06$
- West:  $-1/25$
- Portanto, célula (3,1) assume o valor 0.76 (East)



SCE-5868 Denis F. Wolf

## Resultado após convergência

1	0.812 ↑	0.868 →	0.912 →	Reward +1
2	0.762 ↑	Blocked CELL	0.660 ↑	Reward -1
3	0.705 ↑	0.655 ←	0.611 ↑	0.388 ←
	1	2	3	4

SCE-5868 Denis F. Wolf

### Resultado após convergência

1				Reward +1
2				Reward -1
3				
	START	2	3	4

SCE-5868 Denis F. Wolf

### Value Iteration for Motion Planning

Exemplo de plano obtido utilizando Value Iteration

SCE-5868 Denis F. Wolf

### Value Iteration for Motion Planning

SCE-5868 Denis F. Wolf

## Processo de Decisão de Markov Parcialmente Observável - PDMPO

SCE-5868 Denis F. Wolf

### PDMPO

- ❑ O princípio básico é o mesmo do PDM, porém o estado não é diretamente observável
- ❑ As decisões devem ser tomadas com base em uma distribuição probabilística (*belief* -  $b$ ) do estado
- ❑ PDMPO:
 
$$V_T(b) = \gamma \max_u \left[ r(b, u) + \int V_{T-1}(b') p(b' | u, b) db' \right]$$

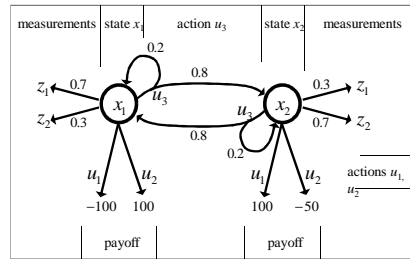
SCE-5868 Denis F. Wolf

### PDMPO

- ❑ PDMPO é uma função de uma distribuição probabilística.
- ❑ PDMPO é um problema complexo, uma vez que distribuições probabilísticas são funções contínuas
- ❑ Para cenários finitos, com um número finito de ações, estados e horizontes, é possível representar as funções de valores (value functions) por segmentos de funções lineares.

SCE-5868 Denis F. Wolf

## PDMPO - Exemplo



SCE-5868 Denis F. Wolf

25

## Parâmetros do exemplo

- As ações  $u_1$  e  $u_2$  são ações terminais.
- A ação  $u_3$  é uma ação de percepção que potencialmente leva a uma mudança de estados.
- O horizonte é finito e  $\gamma=1$ .

$$\begin{aligned} r(x_1, u_1) &= -100 & r(x_2, u_1) &= +100 \\ r(x_1, u_2) &= +100 & r(x_2, u_2) &= -50 \\ r(x_1, u_3) &= -1 & r(x_2, u_3) &= -1 \end{aligned}$$

$$\begin{aligned} p(x'_1|x_1, u_3) &= 0.2 & p(x'_2|x_1, u_3) &= 0.8 \\ p(x'_1|x_2, u_3) &= 0.8 & p(x'_2|x_2, u_3) &= 0.2 \end{aligned}$$

$$\begin{aligned} p(z_1|x_1) &= 0.7 & p(z_2|x_1) &= 0.3 \\ p(z_1|x_2) &= 0.3 & p(z_2|x_2) &= 0.7 \end{aligned}$$

SCE-5868 Denis F. Wolf

## Payoff in POMDPs

- Em um PDM, a recompensa depende do estado do sistema.
- Em um PDMPO, o estado do sistema não é conhecido.
- A recompensa esperada é calculada através da integração de todos os estados:

$$\begin{aligned} r(b, u) &= E_x[r(x, u)] \\ &= \int r(x, u)p(x) dx \\ &= p_1 r(x_1, u) + p_2 r(x_2, u) \end{aligned}$$

SCE-5868 Denis F. Wolf

## Recompensa - exemplo

- Se há 100% de certeza que o estado é  $x_1$  então, quando executada a ação  $u_1$ , recebemos a recompensa de -100
- Se há 100% de certeza que o estado é  $x_2$  então, quando executada a ação  $u_1$ , recebemos a recompensa de +100.
- Para valores intermediários, é calculada uma função linear:

$$\begin{aligned} r(b, u_1) &= -100 p_1 + 100 p_2 \\ &= -100 p_1 + 100 (1 - p_1) \end{aligned}$$

$$r(b, u_2) = 100 p_1 - 50 (1 - p_1)$$

$$r(b, u_3) = -1$$

SCE-5868 Denis F. Wolf

## Função resultante

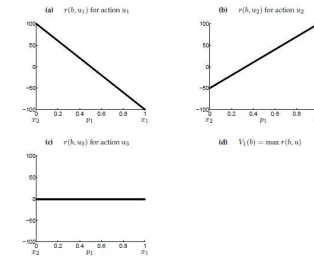
- A value function resultante  $V_1(b)$  corresponde ao máximo das 3 funções a cada ponto.

$$\begin{aligned} V_1(b) &= \max_u r(b, u) \\ &= \max \left\{ \begin{array}{l} -100 p_1 + 100 (1 - p_1) \\ 100 p_1 - 50 (1 - p_1) \\ -1 \end{array} \right\} \end{aligned}$$

- A função linear é composta por segmentos de retas em um formato convexo.

SCE-5868 Denis F. Wolf

## Recompensa - exemplo



SCE-5868 Denis F. Wolf

30

## Plano para T=1

- Dado que o PDMPPO tem horizonte finito T=1, utiliza-se  $V_1(b)$  para determinar o melhor plano.
- No exemplo utilizado, o melhor plano para T=1 é:

$$\pi_1(b) = \begin{cases} u_1 & \text{if } p_1 \leq \frac{3}{7} \\ u_2 & \text{if } p_1 > \frac{3}{7} \end{cases}$$

SCE-5868 Denis F. Wolf

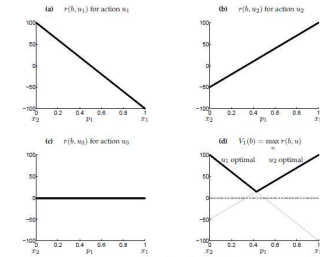
## Pruning

- Ao analisar  $V_1(b)$ , pode-se notar que apenas os 2 primeiros componentes contribuem para a solução final.
- O terceiro componente pode ser eliminado (pruned) de  $V_1(b)$  sem interferir no resultado final.

$$V_1(b) = \max \left\{ \begin{matrix} -100 p_1 & +100 (1 - p_1) \\ 100 p_1 & -50 (1 - p_1) \end{matrix} \right\}$$

SCE-5868 Denis F. Wolf

## Recompensa - exemplo



SCE-5868 Denis F. Wolf

33

## Integrando a percepção

- Considerar que o robô pode fazer uma observação antes de decidir pela ação.
- $p'_i = p(x_i | z) = p(z | x_i) \cdot p(x_i) / p(z)$
- Supondo que o robô observa  $z_i$  para o qual  $p(z_i | x_1) = 0.7$  e  $p(z_i | x_2) = 0.3$ .
- Dada a observação  $z_i$ , o sistema é atualizado usando a regra de Bayes].

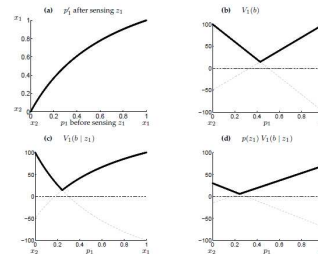
$$p'_1 = \frac{0.7 p_1}{p(z_i)}$$

$$p'_2 = \frac{0.3(1-p_1)}{p(z_i)}$$

$$p(z_i) = 0.7 p_1 + 0.3(1-p_1) = 0.4 p_1 + 0.3$$

SCE-5868 Denis F. Wolf

## PDMPPO - Exemplo



SCE-5868 Denis F. Wolf

35

## Valor esperado depois da percepção

- Como não se sabe qual será a próxima observação, deve-se computar o belief esperado:

$$\begin{aligned} \bar{V}_1(b) &= E_z[V_1(b | z)] \\ &= \sum_{i=1}^2 p(z_i) V_1(b | z_i) \\ &= \max \left\{ \begin{matrix} -70 p_1 & +30 (1 - p_1) \\ 70 p_1 & -15 (1 - p_1) \end{matrix} \right\} \\ &\quad + \max \left\{ \begin{matrix} -30 p_1 & +70 (1 - p_1) \\ 30 p_1 & -35 (1 - p_1) \end{matrix} \right\} \end{aligned}$$

SCE-5868 Denis F. Wolf

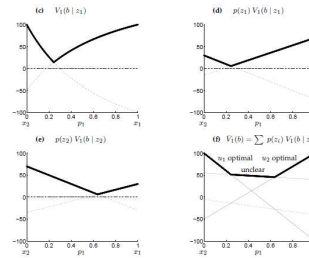
## Valor esperado depois da percepção

- As 4 possibilidades consistem em:

$$\begin{aligned} \bar{V}_1(b) &= \max \begin{Bmatrix} -70 p_1 + 30(1-p_1) & -30 p_1 + 70(1-p_1) \\ -70 p_1 + 30(1-p_1) & +30 p_1 - 35(1-p_1) \\ +70 p_1 - 15(1-p_1) & -30 p_1 + 70(1-p_1) \\ +70 p_1 - 15(1-p_1) & +30 p_1 - 35(1-p_1) \end{Bmatrix} \\ &= \max \begin{Bmatrix} -100 p_1 + 100(1-p_1) \\ +40 p_1 + 55(1-p_1) \\ +100 p_1 - 50(1-p_1) \end{Bmatrix} \end{aligned}$$

SCE-5868 Denis F. Wolf

## PDMPO - Exemplo



SCE-5868 Denis F. Wolf

38

## Transição de estados

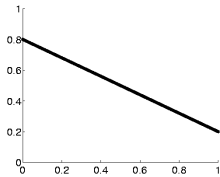
- Quando a ação  $u_3$  o estado potencialmente muda, o que deve ser levado em consideração na predição.

$$\begin{aligned} p'_1 &= E_x[p(x_1 | x, u_3)] \\ &= \sum_{i=1}^2 p(x_1 | x_i, u_3) p_i \\ &= 0.2 p_1 + 0.8(1-p_1) \\ &= 0.8 - 0.6 p_1 \end{aligned}$$

SCE-5868 Denis F. Wolf

## Transição de estados

$$\begin{aligned} p'_1 &= E_x[p(x_1 | x, u_3)] \\ &= \sum_{i=1}^2 p(x_1 | x_i, u_3) p_i \\ &= 0.2 p_1 + 0.8(1-p_1) \\ &= 0.8 - 0.6 p_1 \end{aligned}$$



SCE-5868 Denis F. Wolf

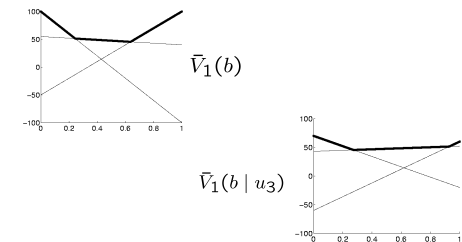
## Value Function depois de executar $u_3$

- Considerando a possível transição de estados obtemos:

$$\begin{aligned} \bar{V}_1(b) &= \max \begin{Bmatrix} -70 p_1 + 30(1-p_1) & -30 p_1 + 70(1-p_1) \\ -70 p_1 + 30(1-p_1) & +30 p_1 - 35(1-p_1) \\ +70 p_1 - 15(1-p_1) & -30 p_1 + 70(1-p_1) \\ +70 p_1 - 15(1-p_1) & +30 p_1 - 35(1-p_1) \end{Bmatrix} \\ &= \max \begin{Bmatrix} -100 p_1 + 100(1-p_1) \\ +40 p_1 + 55(1-p_1) \\ +100 p_1 - 50(1-p_1) \end{Bmatrix} \\ \bar{V}_1(b | u_3) &= \max \begin{Bmatrix} 60 p_1 - 60(1-p_1) \\ 52 p_1 + 43(1-p_1) \\ -20 p_1 + 70(1-p_1) \end{Bmatrix} \end{aligned}$$

SCE-5868 Denis F. Wolf

## Value Function depois de executar $u_3$



SCE-5868 Denis F. Wolf



## Porque Pruning é essencial

- Cada passo adiciona novos componentes lineares para  $V$ .
- Cada observação eleva ao quadrado o número de componentes lineares.
- Exemplo: Uma função "un-pruned" com  $T=20$  possui mais de  $10^{547,864}$  componentes lineares.
- Com  $T=30$ , temos  $10^{561,012,337}$  componentes lineares.
- A função pruned com  $T=20$ , contém somente 12 componentes lineares.
- A explosão combinatória de componentes lineares torna os PDMPPO impraticáveis em diversas aplicações.

SCE-5868 Denis F. Wolf

## PDMPPO - Algoritmo

```

1: Algorithm PDMPPO:
2:  $\mathcal{T} = \{0, \dots, 0\}$ 
3: for  $\tau = 1$  to  $T$  do
4:    $\mathcal{T} = \emptyset$ 
5:   for all  $(s^1; \dots; s^N)$  in  $\mathcal{T}$  do
6:     for all control actions  $u$  do
7:       for all measurements  $z$  do
8:         for  $j = 1$  to  $N$  do
9:            $v_{j, \tau}^u = \sum_{s^j} p^j(s^j | s^j, u, x_j)$ 
10:        endfor
11:      endfor
12:    endfor
13:  endfor
14:  for all control actions  $u$  do
15:    for all  $k(1), \dots, k(M) = (1, \dots, 1)$  to  $(|\mathcal{T}|, \dots, |\mathcal{T}|)$  do
16:      for  $i = 1$  to  $N$  do
17:         $v_i^u = \gamma \left[ v_{k(i), \tau}^u + \sum_{s^i} A_{s^i}^i(v_{k(i), \tau}^u) \right]$ 
18:      endfor
19:      add  $(s^1, \dots, s^N)$  to  $\mathcal{T}$ 
20:    endfor
21:  endfor
22:  optional: prune  $\mathcal{T}$ 
23:   $\mathcal{T} = \mathcal{T}$ 
24: endfor
25: return  $\mathcal{T}$ 
  
```

SCE-5868 Denis F. Wolf

44

## PDMPPO - Sumário

- PDMPPOs computam ações ótimas para ambientes estocásticos e parcialmente observáveis.
- Para problemas de horizonte finito, a função resultante é composta por segmentos de reta convexas..
- A cada iteração, o número de componentes lineares cresce exponencialmente.
- PDMPPOs vem sendo aplicados para situações com poucos estados possíveis e poucas ações e observações possíveis.

SCE-5868 Denis F. Wolf

## Point-based Value Iteration

- Mantém apenas alguns possíveis estados válidos
- Somente considera situações que maximizem as funções de valor de pelo menos um desses estados

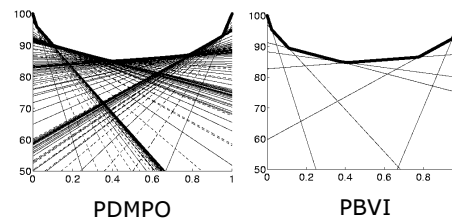
Ex:  $B = \{p_1=0.0, p_1=0.1, p_1=0.2, \dots, p_1=1.0\}$

Para  $T=30$  PDMPPO =  $10^{561,012,337}$  e PBVI = 120

SCE-5868 Denis F. Wolf

## Point-based Value Iteration

Value functions for  $T=30$



SCE-5868 Denis F. Wolf