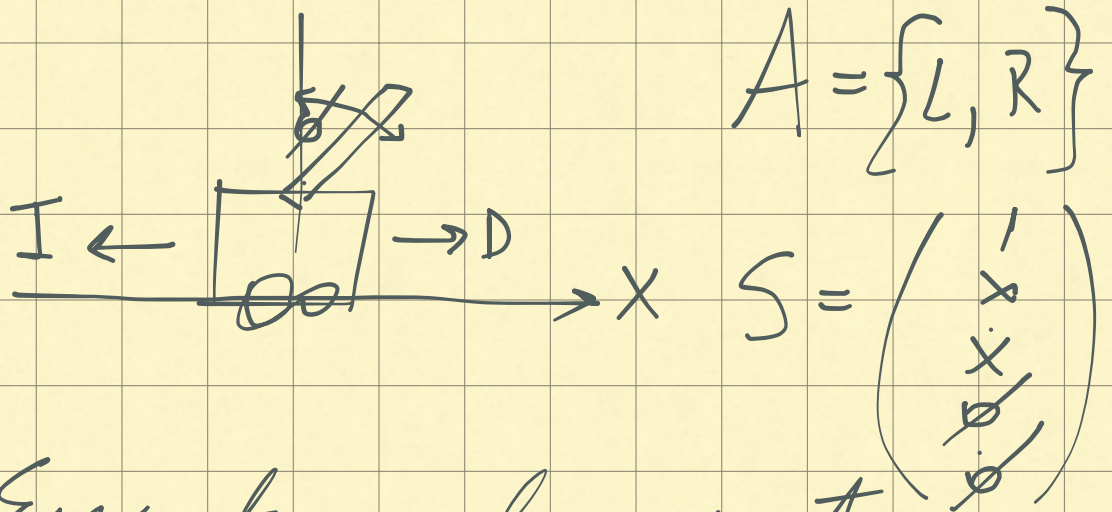


Clase 5.1 "Policy Search", POMDP

$$V^* \rightarrow \pi^*$$

Búsqueda de Política Directa

Caso de uso "Cart Pole"



Paso 1: Encontrar el conjunto de funciones para obtener una π aproximada.

Reg. log.

$$\textcircled{y} \quad \underline{h_\theta(x)} = \frac{1}{1 + e^{-\theta^T x}}$$

Buscar la π directamente

$\mathcal{R} \sim \pi_{\theta}(s)$ — depende de un conjunto de estados
 dado un conj. de parámetros

$$\pi_{\theta}(s) = \frac{1}{1 + e^{-\theta^T s}}$$

$$y \quad s = \begin{pmatrix} 1 \\ x \\ x \\ \cancel{\phi} \\ \cancel{\phi} \end{pmatrix}$$

Paso: Elegir una Política (Estocástica)

función $\pi: S \times A \rightarrow \mathbb{R}$

$\pi(s, a)$

Ejemplo

$$\pi_{\theta}(s, "D") = \frac{1}{1 + e^{-\theta^T s}}$$

$$\pi_{\theta}(S, "I") = 1 - \frac{1}{1 + e^{-\theta^T S}}$$

$$\epsilon = 0.3$$

ϵ -greedy

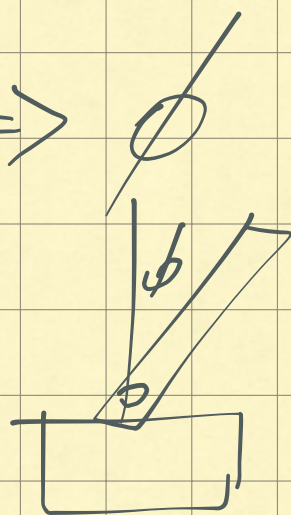
30% Aleatorio

70% greedy

$$S = \begin{pmatrix} 1 \\ x \\ x \\ \phi \\ \phi \end{pmatrix}$$

$$\theta = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

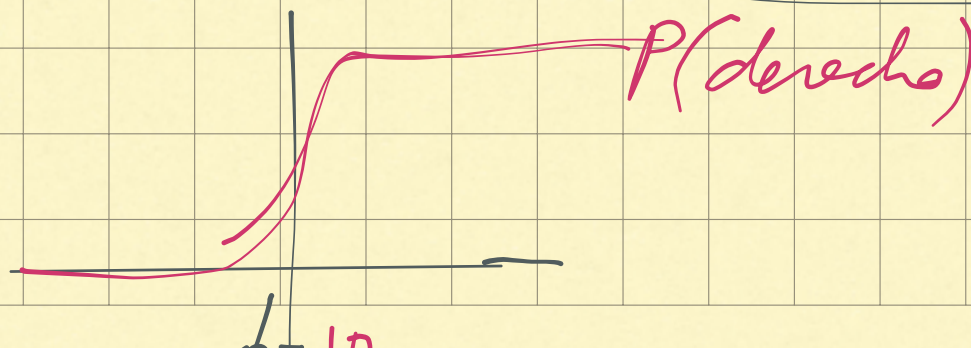
\Rightarrow



En este caso, $\theta^T = \phi$

$$\pi_{\theta}(S, "D") = \frac{1}{1 + e^{-\phi}}$$

CartPole



Un poco más complejo:

\times
 ϕ

$$\theta = \begin{bmatrix} 0 \\ -0.5 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Cuanto debe
acelerar a la
derecha,
dependiendo de
estos 2 parámetros

(pos, vel. lateral, ángulo, vel. angular)

Policy Search

- El Objetivo es encontrar este conjunto de parámetros θ
- Para ejecutar $\pi_\theta(s, a)$
- Maximicemos la ganancia

$$\max_{\theta} E[R(s_0, a_0) + \dots + R(s_t, a_t) | \pi_\theta]$$

Encontrar θ :

→ Horizonte Finito T

→ Estado Inicial S_0

$$\max_{\theta} E \left[R(S_0, a_0) + R(S_1, a_1) + \dots + R(S_t, a_t) \mid T \right]$$

$T=1$

$$= \sum P(S_0, a_0, S_1, a_1) \left[\overbrace{R(S_0, a_0) + R(S_1, a_1)}^{\text{ganancia}} \right]$$

$$= \sum P(S_0) \pi_{\theta}(S_0, a_0) \underbrace{P(S_1)}_{S_0, a_0} \pi_{\theta}(S_1, a_1) [g]$$

Derivar en fun de θ ,
usando algo. de gradientes
ascendentes (estocásticos)

θ/θ

~~CONFUSION~~

Algoritmo Reinforce (Reinforce Algorithm)

→ Monte Carlo
Policy Gradient

{
sample $S_0, A_0, S_1, A_1, \dots$ $\xrightarrow{T=1}$
compute $R(S_0) + R(S_1)$
update

$$\Theta := \Theta + \alpha \left[\frac{\nabla_{\Theta} \Pi_{\Theta}(S_0, A_0)}{\Pi_{\Theta}(S_0, A_0)} + \frac{\nabla_{\Theta} \Pi_{\Theta}(S_1, A_1)}{\Pi_{\Theta}(S_1, A_1)} \right] g$$

(Hacia el aprendizaje)

}

En cada iteración, actualizamos Θ

Entonces, El Algoritmo Reinforce

- las actualizaciones tienen aleatoriedad

- depende de la secuencia de estado

- Hasta que logre el Target

$$\nabla_{\theta} \left[\sum_{S, a, S'} P(S_0) \pi_{\theta}(S_0, a_0) P_{S_0}(S_1) \pi_{\theta}(S_1, a_1) (g) \right]$$

Aplicamos "Regla del Producto"

$$\frac{d}{d\theta} f(\theta) g(\theta) h(\theta)$$

$$= f'(\theta) g(\theta) h(\theta) + f(\theta) g'(\theta) h(\theta) + f(\theta) g(\theta) h'(\theta)$$

$$\sum \left[P(S_0) \pi_{\theta}(S_0, a_0) \nabla_{\theta} \pi_{\theta}(S_0, a_0) P(S_1) \pi_{\theta}(S_1, a_1) \right]$$

$$\pi_\theta(s_0, a_0)$$

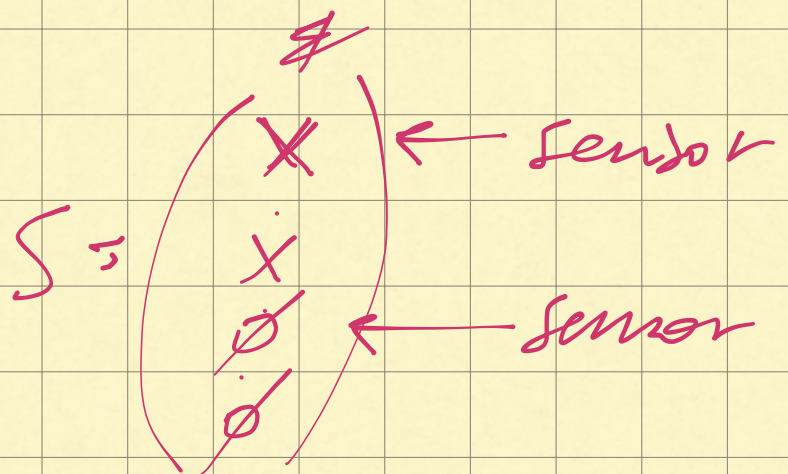
$$+ \left[P(s_0) \pi_\theta(s_0, a_0) P(s_1) \pi_\theta(s_1, a_1) \frac{\nabla_\theta \pi_\theta(s_1, a_1)}{\pi_\theta(s_1, a_1)} \right]$$

factorizando

$$= \sum P(s_0, a_0, s_1, a_1) \left[\underbrace{\frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \frac{\nabla_\theta \pi_\theta(s_1, a_1)}{\pi_\theta(s_1, a_1)}}_{\text{Gradient Update}} \right] g$$

Conclusions:

POMDP



- Siempre tengo una medición/observación parcial, ruidosa del estado

2/4

- Generalization

- Función muy bien POMDP

$$y = \begin{pmatrix} x \\ \phi \end{pmatrix} + \text{ruido}$$

$$\pi_{\theta}(y, "D") = \frac{1}{1 + e^{-\theta^T y}}$$

