

Agenda

1. Depuración de los algoritmos de RL
2. Algoritmos de Búsqueda de Política Directa (Policy Search)
3. *Multi-arm Bandits*
4. Deep Q Learning
5. Taller práctico (Tarea incluida)

PPT

Guaderno

Implementación

Subir Solución

IT Value

MA Bandit

Q-learning

Clase 6

Métodos Avanzados, Impl., Apps.

Depuración de Algoritmos de RL

Clase 5.1 Jorge Vásquez

Algoritmo de ML

basados en \rightarrow Data \rightarrow functions \leftarrow lineales parametrizados $f_{\theta} \rightarrow NN$

Model



Q values
 $R(s_0, a_0) + \gamma \dots$
 $R(a_0) + \gamma R(a_{\dots})$

1. Construir un modelo / simulador de un helicóptero
2. Escoger una función de recompensas, por ej. $R(s) = -|S - S_{desired}|^2$
3. Hacer correr un algoritmo de RL para volver el helicóptero en un simulador y maximizar $E[R(\underline{S_0}) + R(S_1) + \dots + R(S_t)]$ y luego obtener una política π

$(a_0) + \dots$
 $(s, a) + \dots$

Algoritmo de ML



74

1. Supongamos que hacemos eso y el controlador resultante π entrega un rendimiento mucho peor que un piloto humano. ¿Qué hacemos?

✓ ¿Mejorar el modelo o simulador?

• ¿Modificar la función de recompensas R ?

✓ ¿Modificar el algoritmo de RL?

Estática $R_0 + R_1 + \dots + R_n$
Dinámica

$\max E(\quad)$ ✓
 $- A, S, \delta, H,$

Depurar un Algoritmo de RL

- El controlador dado π rinde mal ✓
- Asumamos : MODEL
 - 1 • El simulador del helicóptero es preciso ✓
 - 2 • El algoritmo de RL controlar correctamente el helicóptero en el simulador y también maximiza ganancia esperada
 - payoff $V^{\pi_R}(s_0) = E[R(s_0) + R(s_1) + \dots + R(s_T) \mid \pi_{RL}, s_0]$. ✓
 - 3 • Maximizar ganancia esperada corresponde a un vuelo autónomo correcto. ✓
- Luego, el controlador π debería hacer volar bien el helo

Depurar un Algoritmo de RL

- Diagnosticar:

1. Si la π_{RL} vuela bien en el simulador, pero no en el mundo real, entonces el problema es el simulador.

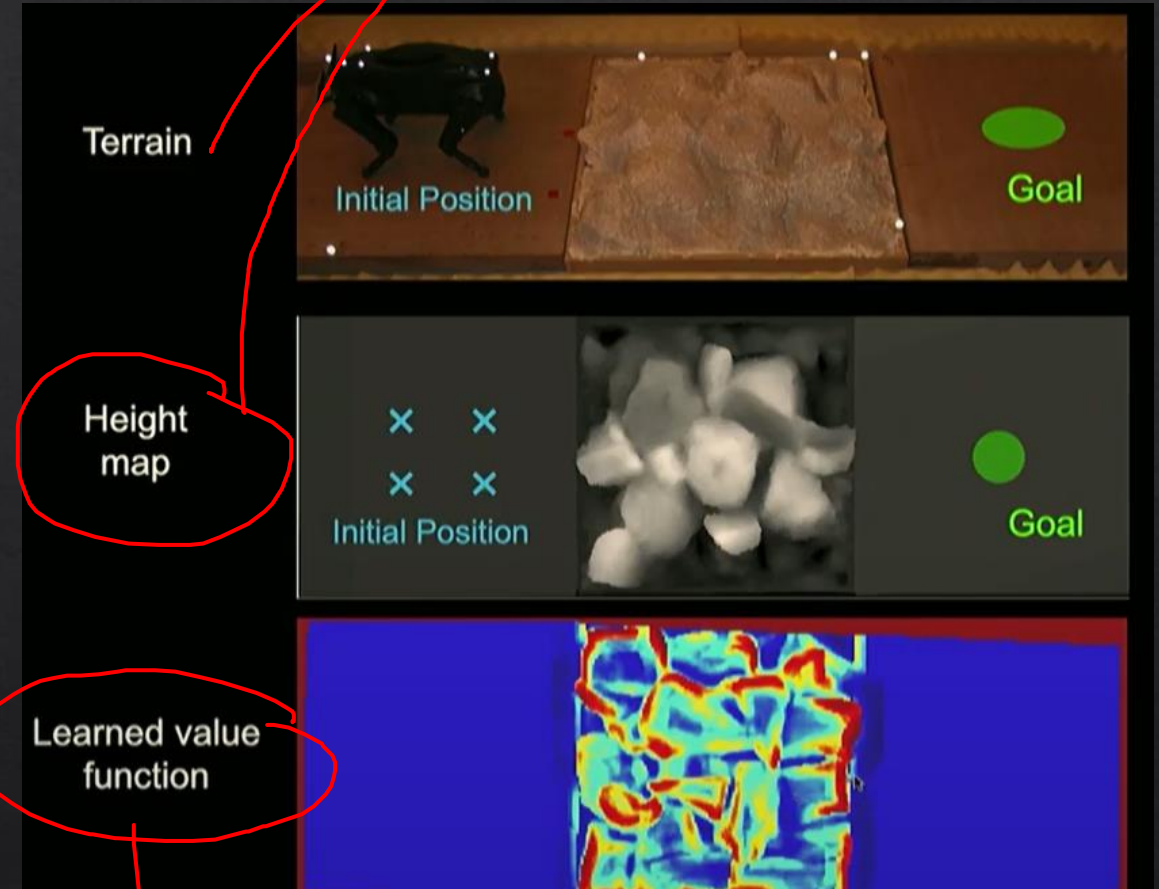
no delo

$$V^{\pi_{RL}}(s_0) < V^{\pi_{humano}}(s_0)$$

2. Dejemos que π_{humano} sea la política de control y si entonces el problema es el algoritmo de RL (no maximiza la ganancia esperada)

3. De lo contrario, el problema es en la función de costo (la maximización no corresponde a un buen vuelo autónomo)

Perro Robótico



II. fn. de pérdida

Algo RL \rightarrow evaluar