

MDPs: overview

Son un modelo matemático para la toma de decisiones bajo incertidumbre, este se desarrollo en los años 50-60 y están relacionados con las cadenas de Markov y permiten tomar decisiones en sistemas estocásticos.

Sus aplicaciones son:

- Robótica: decidir movimientos con actuadores defectuosos o obstáculos inesperados.
- Asignación de Recursos: decidir qué producir sin conocer la demanda exacta.
- Agricultura: decidir qué plantar sin conocer el clima futuro.

Definición de un MDP:

- **Estados (S):** conjunto de posibles estados del sistema.
- **Acciones (A):** conjunto de acciones disponibles en cada estado.
- **Transición (T(s, a, s')):** probabilidad de pasar de estado sss a estado s's's' al tomar acción aaa.
- **Recompensa (R(s, a, s')):** recompensa obtenida en una transición.
- **Factor de descuento (γ gamma):** ponderación de recompensas futuras (valor entre 0 y 1).

Evaluación de Políticas

- **Política (π pi):** función que asigna una acción a cada estado.
- **Valor de una política $V_{\pi}(s)$ $V_{\pi}(s)$:** utilidad esperada de seguir la política desde sss.
- **Q-valor $Q_{\pi}(s,a)$ $Q_{\pi}(s,a)$:** utilidad esperada de tomar acción aaa en sss y seguir la política después.
- Se puede calcular mediante un algoritmo iterativo hasta la convergencia.

Algoritmos para Resolver MDPs

Evaluación de Políticas

- Se inicializa $V(s)=0$ $V(s)=0$
- Se actualizan los valores usando:

$$V_{\pi}(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_{\pi}(s')] \\ V_{\pi}(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_{\pi}(s')]$$

- Se repite hasta convergencia.

Iteración de Valor (Value Iteration)

- Se busca la mejor política optimizando el valor esperado:

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \quad V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

- Se actualizan los valores de los estados iterativamente.

Iteración de Política

- Se alterna entre evaluar una política y mejorarla iterativamente.