

## Inteligencia Artificial

### Tema 4 – Procesos de Decisión de Markov

#### 1. Introducción

Vamos a tratar secuencias de acciones cuyos efectos son inciertos: similares a espacios de estados, pero el efecto de una acción está descrito como una distribución de probabilidad. El resultado de una acción es determinista. Además, se introduce el concepto de “recompensa” para un estado.

Se busca cual es la mejor acción a realizar en cada momento, ¿cómo se decide? → sabiendo de antemano sus posibles efectos y con que probabilidad, además de la recompensa en cada caso.

#### 2. El modelo matemático

Un proceso de decisión de Markov viene definido por:

- a. Un conjunto  $S$  de estados (con un estado inicial  $S_0$ )
- b. Para cada estado, un conjunto  $A(s)$  de acciones aplicables al estado
- c. Un modelo de transición, dado por una distribución de probabilidad  $P(s' / s, a)$  para cada par de estados  $s', s$  y acción aplicable a  $s$   $a$ , que indica la probabilidad de que aplicando  $a$  a  $s$  se obtenga  $s'$ .
- d. Una función de recompensa  $R(s)$ .

Propiedad de Markov: el efecto (incierto) de una acción sobre un estado solo depende de la acción y de propio estado (no de estados anteriores). En consecuencia, como hemos dicho, una solución nunca puede ser una secuencia de acciones ya que efecto de cada una es incierto. Lo que buscamos más bien es una “política” que en cada posible estado recomiende que acción aplicar: formalmente una política es una función  $\pi$  definida sobre el conjunto de estados  $S$  tal que  $\pi(s) \in A(s)$ .

Una política puede generar secuencias de acciones distintas (unas con más probabilidad que otras) pero buscamos siempre la política “óptima” que maximice la recompensa media esperada para las posibles secuencias de acciones que se generan.

Valoración de una secuencia de estados en el tiempo: suponiendo que mediante una aplicación se genera una secuencia de estados  $q_0, q_1, q_2, \dots$ . ¿Cómo valoramos una secuencia de estados? Usamos la valoración mediante recompensa con descuento:

$$V([q_0, q_1, q_2, \dots]) = R(q_0) + \gamma * R(q_1) + \gamma^2 * R(q_2) + \dots$$

donde  $\gamma$  es el llamado “factor de descuento”.

Las valoraciones en general no son infinitas en tanto que, si hay estados terminales, los asimilamos a estados a partir de los cuales las recompensas son cero. Aún así, si las recompensas están acotadas por una cantidad  $R_{max}$  e  $\gamma < 1$ , entonces la valoración de una secuencia no puede ser mayor a  $R_{max}/(1 - \gamma)$ .

Dado una política  $\pi$  y un estado  $s$ , podemos valorar  $s$  respecto de  $\pi$  teniendo en cuenta la valoración de las secuencias de estados que se generan si se sigue dicha política a partir de  $s$ . La valoración de un estado respecto de una política  $\pi$  la denotamos como  $V^\pi(s)$ .

Por suerte, para calcular  $V^\pi(s)$  no necesitamos calcular todas las posibles secuencias que podrían iniciarse en  $s$ . Calcular  $V^\pi$  es resolver el siguiente sistema de ecuaciones lineales:

$$V^\pi(s) = R(s) + \gamma \cdot \sum_{s'} (P(s'|s, \pi(s)) \cdot V^\pi(s'))$$

donde las incógnitas son los  $V^\pi$ , una por cada estado  $s$ .

Llamamos muestreo a generar secuencias de estados aplicando las acciones que indica la política teniendo en cuenta la distribución de probabilidad y calcular la media de sus valoraciones:

- Comenzamos por un valor arbitrario  $V_0^\pi(s)$  para cada estado  $s$ .
- En cada iteración se aplica la fórmula que relaciona la valoración de un estado con la de sus vecinos:

$$V_{i+1}^\pi(s) = R(s) + \gamma \cdot \sum_{s'} (P(s'|s, \pi(s)) \cdot V_i^\pi(s'))$$

Políticas óptimas y valoraciones de estados:

- La valoración de un estado  $s$  denotada como  $V(s)$  se define como:  $V(s) = \max(\pi) V^\pi(s)$
- La política óptima se define como  $\pi^*$  tal que  $\pi^*(s) = \operatorname{argmax}(\sum_{a \in A(s)} P(s'|s, a) \cdot V(s'))$  donde  $a$  pertenece a  $A(s)$
- $V^{\pi^*} = V$  y nuestro objetivo es calcular  $V$  y  $\pi^*$

### 3. Ecuaciones de Bellman

De manera análoga a  $V^\pi$ , describimos  $V(s)$  en función de las valoraciones de los estados vecinos. Se conoce esto como ecuaciones de Bellman:

$$V(s) = R(s) + \gamma \cdot \max_{a \in A(s)} \sum_{s'} (P(s'|s, a) \cdot V(s'))$$

La solución a este sistema (una ec. por estado) nos da la valoración de cada estado, y en consecuencia, la política óptima. Para resolver las ecuaciones de Bellman planteamos un método iterativo, el muestreo visto arriba. Se puede demostrar que los valores que se obtienen convergen asintóticamente hacia la solución (única si  $\gamma < 1$ ) de las ecuaciones. La convergencia es rápida para valores de  $\gamma$  pequeños y viceversa.

Algoritmo de iteración de valores para calcular las ecuaciones de Bellman:

### Entrada:

- Un proceso de decisión de Markov: conjunto de estados  $\mathbf{S}$ ,  $\mathbf{A}(\mathbf{s})$ , modelo de transición  $\mathbf{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ , recompensa  $\mathbf{R}(\mathbf{s})$  y descuento  $\gamma$
- $\epsilon > 0$ , cota de error máximo permitido

### Salida:

- Una valoración  $\mathbf{V}(\mathbf{s})$  para cada estado

### Procedimiento:

- Inicio: Sea  $\mathbf{V}_0$  una función sobre los estados, con valor 0 para cada estado,  $\delta = \infty$  e  $i$  igual a 0
- Repetir
  - Hacer  $i = i + 1$  y  $\delta = 0$
  - Para cada  $\mathbf{s} \in \mathbf{S}$  hacer:
    - $\mathbf{V}_i(\mathbf{s}) \leftarrow \mathbf{R}(\mathbf{s}) + \gamma \cdot \max_{a \in \mathbf{A}(\mathbf{s})} \sum_{\mathbf{s}'} (\mathbf{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) \cdot \mathbf{V}_{i-1}(\mathbf{s}'))$
    - Si  $|\mathbf{V}_i(\mathbf{s}) - \mathbf{V}_{i-1}(\mathbf{s})| > \delta$  entonces  $\delta = |\mathbf{V}_i(\mathbf{s}) - \mathbf{V}_{i-1}(\mathbf{s})|$
- Devolver  $\mathbf{V}_i$

#### 4. Iteración de Políticas

Existe un método alternativo de obtener  $V$ , comenzando con una política inicial  $\pi_0$ :

- Calcular  $V^{\pi_i}$ .
- A partir de  $V^{\pi_i}$  calcular una nueva política  $\pi_{i+1}$  que en cada estado recomiende la acción que mejor valoración espera, respecto de  $V^{\pi_i}$ .
- Los dos pasos anteriores se iteran hasta que  $\pi_i$  y  $\pi_{i+1}$  coinciden, en cuyo caso hemos alcanzado la política óptima y su valoración asociada, es  $V$ .

El proceso anterior termina en tanto que hay un nº finito de políticas distintas. El algoritmo de iteración de políticas:

### Entrada:

- Un proceso de decisión de Markov: conjunto de estados  $\mathbf{S}$ ,  $\mathbf{A}(\mathbf{s})$ , modelo de transición  $\mathbf{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ , recompensa  $\mathbf{R}(\mathbf{s})$  y descuento  $\gamma$
- Un número  $k$  de iteraciones

**Salida:** Una política óptima  $\pi$  y su valoración asociada

### Procedimiento:

- Inicio:  $\pi$  una política aleatoria, asignando una acción a cada estado
- Repetir

- Calcular  $V^\pi$

- Hacer **actualizada** igual a **Falso**

- Para cada estado  $\mathbf{s}$ :

- Si

$$\max_{a \in A(s)} \sum_{s'} (P(s' | s, a) \cdot V^\pi(s')) > \sum_{s'} (P(s' | s, \pi(s)) \cdot V^\pi(s')),$$

entonces hacer  $\pi(\mathbf{s})$  igual a  $\mathbf{argmax}_{a \in A(s)} \sum_{s'} (P(s' | s, a) \cdot V^\pi(s'))$  y

**actualizada** igual **Verdad**

hasta que **actualizada** sea **Falso**

- Devolver  $\pi$  y  $V^\pi$