

Capítulo 4: Dynamic Programming.

Exercise 4.1: If π is the equiprob. random policy, what is $q_{\pi}(11, \downarrow)$?

And $q_{\pi}(7, \downarrow)$

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')] = \text{action}$$

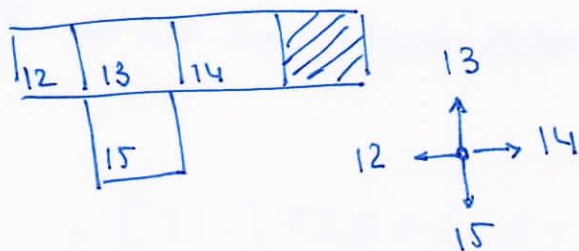
$$q_{\pi}(11, \downarrow) = p(\text{terminal}, -1 | 11, \downarrow) [-1 + \gamma V_{\pi}(\text{terminal})] = -1$$

"
"
"
1.0
1.0
0.0

$$q_{\pi}(7, \downarrow) = p(11, -1 | 7, \downarrow) [-1 + \gamma V_{\pi}(11)] = -15$$

"
"
"
1.0
1.0
-14

Exercise 4.2. We add a state "15" below "13"



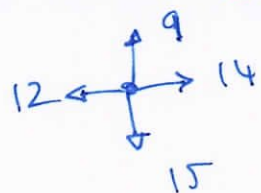
$V_{\pi}(15)$ para $\pi(a|s) = 0.25 \forall s, \forall a \in A(s)$?

$$V_{\pi}(15) = \sum_a \pi(a|s) \sum_{s', r} p(r, s' | s, a) [r + \gamma V_{\pi}(s')] =$$

$$= 0.25 [-1 + \overset{-22}{V_{\pi}(12)}] + 0.25 [-1 + \overset{-20}{V_{\pi}(13)}] + 0.25 [-1 + \overset{-14}{V_{\pi}(14)}] + 0.25 [-1 + V_{\pi}(15)] \neq \dots$$

→ Hay que resolver la ecuación, y queda $V_{\pi}(15) = -20$

b). Supongamos ahora que las distancias en ~~el~~ $s=13$ cambian



¿ $V_n(15)$ ahora?

→ En este caso, $V_n(13)$ cambia! y por ~~tanto~~, $V_n(15)$ también!
Empezaré con $V_n(15) = -20$, ~~y después~~ mejor $V_n(13)$ y luego
revelaré $V_n(15)$? o lo escribo todo en términos de $V_n(15)$

$$V_n(13) = 0.25 [-1 + \overset{V_n(12) = -22}{\cancel{(-22)}}] + 0.25 [-1 + \overset{=-20}{V_n(9)}] + 0.25 \cdot$$

$$[-1 + \underset{-14}{V_n(14)}] + 0.25 [-1 + V_n(15)] \quad (*)$$

$$V_n(15) = 0.25 [-1 + V_n(13)] + 0.25 [-1 + \overset{-22}{V_n(12)}] +$$

$$+ 0.25 [-1 + \underset{-14}{V_n(14)}] + 0.25 [-1 + \overset{\text{per}}{V_n(15)}] \quad (**)$$

Sustituyendo (*) en (**), obtendríamos una ec. de 2° grado:

$$V_n(13) = 0.25 \cdot (-22) + 0.25 [-21] + 0.25 [-15] +$$

$$0.25 [-1 + V_n(15)] = -5.75 + (-5.25) + (-3.75) +$$

$$+ 0.25 [-1 + V_n(15)] =$$

$$= -15 + 0.25 V_n(15)$$

$$V_n(15) = 0.25 [\cancel{(-16)} -16 + 0.25 V_n(15)] + (-5.75) +$$

$$+ (\cancel{(-3.75)} -3.75) + 0.25 [-1 + V_n(15)] =$$

$$= -4 + 0.0625 V_n(15) - 11 - 0.25 + 0.25 V_n(15) =$$

$$-13.75$$

$$- \textcircled{1.75}$$

$$= \textcircled{-15.5} + 0.3125 V_n(15)$$

$$0.6875 V_n(15) = \overset{-13.75}{\textcircled{-15.5}} \Rightarrow V_n(15) = \overset{= 20}{\textcircled{20}}$$

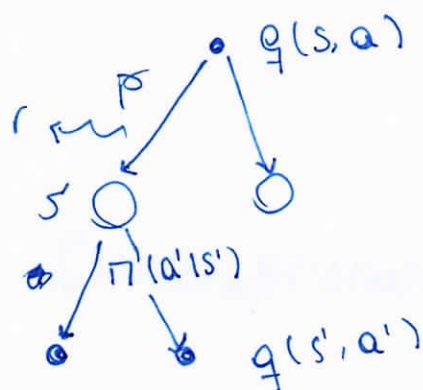
Ex. 4.3 Ecs. análogas a 4.3.4.4 y 4.5 para $q_n(s, a)$?

$$q_n(s, a) = \mathbb{E}_n [G_t | S_t = s, A_t = a] = \mathbb{E}_n \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] =$$

$$= \mathbb{E}_n \left[R_t + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] =$$

$$= \mathbb{E}_n [R_t | S_t = s, A_t = a] + \gamma \mathbb{E}_n [G_{t+1} | S_t = s, A_t = a] =$$

$$= \sum_{s', r} \gamma p(s', r | s, a) + \gamma \sum_{s', r} \pi(a' | s') \sum_{s', r} p(s', r | s, a) q_n(s', a')$$



4 entones:

$$q_{k+1}(s, a) = \sum_{s', r} \gamma p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_k(s', a') \right]$$

Ejercicio 4.4. El alg. de policy iteration en la pág. 80 tiene un bug...

~~For~~ Imaginemos varias π que son todas igual de buenas.
if $V_n(s) \neq V_{n+1}(s)$ then policy-estable \leftarrow false

Ejercicio 4.5. ¿Cómo se definiría la policy iteration para $q_n(s,a)$?

$$q(s,a) \in \mathbb{R} \text{ y } \pi(s)$$

$$V(s) \in \mathbb{R} \text{ arbitrario}$$

$$\pi(s) \in A(s) \text{ "}$$

$$\} \forall s \in \mathcal{S}.$$

Initialization *

$$q_n(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma V_n(s')] \quad \forall s, \forall a$$

Policy Evaluation:

Loop:

$$\Delta \leftarrow 0$$

loop for each $s \in \mathcal{S}$, $\forall a$:

$$q \leftarrow q_n(s,a)$$

$$q_{n+1}(s,a) \leftarrow \sum_{s',r} p(s',r|s,a) [r + \gamma \sum_{a'} \pi(a'|s') q_n(s',a')]$$

$$\Delta \leftarrow \max(\Delta, |q - q_{n+1}(s,a)|)$$

until $\Delta < \theta$

3. Policy Improvement:

policy-stable \leftarrow true

$\forall s \in \mathcal{S}$:

~~old-action~~

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \left[\sum_{s',r} p(s',r|s,a) [r + \gamma q(s',a')] \right]$

if $q(s, \text{old-action}) \neq q(s, \pi(s))$ then policy-stable \leftarrow false

If policy-stable, then stop and return $q \approx q^*$ y $\pi \approx \pi_*$ else 2.

* otra posible inicialización:

$\pi \leftarrow$ arbitrario

$Q \leftarrow$ arbitraria $\mathcal{S} \times \mathcal{A}(s) \mapsto \mathbb{R}$

Ej. 4.6: políticas ϵ -soft: la prob. de escoger cada acción en cada estado es $\frac{\epsilon}{|\mathcal{A}(s)|}$. ¿cambios necesarios a introducir en los pasos 3.2 y 1 del algoritmo de iteración de política para cumplir con esto?

~~Paso 3:~~

~~$\forall s \in \mathcal{S}$ y $\forall a$:~~

~~$\pi(a|s) = \frac{\epsilon}{|\mathcal{A}(s)|}$ if $a \neq \arg \max$~~

~~if $\pi(a|s) = \frac{\epsilon}{|\mathcal{A}(s)|}$ if $a \neq \arg \max$~~

~~$\pi_+ = \frac{\epsilon}{|\mathcal{A}(s)|}$~~

En este caso, la política pasaría a ser estocástica.

En el paso 3, tendríamos que darle una prob. $\pi(a|s) = \frac{\epsilon}{|A(s)|}$ a las acciones que no son óptimas y el resto de la masa

$1 - [|A(s)| - 1] \frac{\epsilon}{|A(s)|}$ a la acción óptima.

En el paso 2, tendríamos

$$V(s) \leftarrow \sum_{s', r} p(s', r|s) \sum_a \pi(a|s) \sum_{s'', r'} p(s'', r'|s', a) [r + \gamma V(s'')]$$

En el paso 1, $\pi(s)$ es ahora una distribución, no es una función:

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

Ejercicio 4.8. ¿Por qué la política óptima del problema del jugador viene en forma? ¿~~Qué otras~~ & ...

Estados: capital em cada momento. $S = \{1, 2, \dots, 99\}$

$$a \in \{0, 1, \dots, \min(s, 100-s)\}$$

Ex. $\therefore S = 54$.

$$a \in \{0, 1, \dots, \min(\pi_{4,100} - 54)\} = \{0, 1, \dots, 46\}$$

¿Por qué es mío? Para no ganar más de los:

Tengo 54. Después 54 \rightarrow Gano 54 + 54 = 108.
máximo

En lugar de eso, lo que gano es lo que me falta para llegar a los:

$$100 - 54 = 46.$$

~~Si lo que apuesto es 33, que es menor que 50, entonces gano $33+33=$~~
~~No tengo que~~

Si $s = 33$ y apesto 33, entonces no hay que hacer el add más anterior, ya que $33 + 33 < 100$.

En general, cuando $s > 50$ y lo apuesto todo, lo que gano es

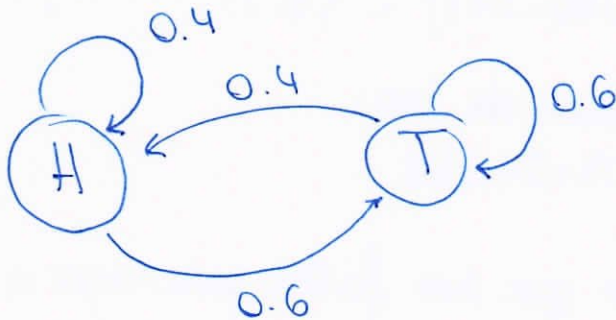
~~Soe~~ $s + (100 - s) = 100$

$$P_h = 0.4 ; p_t = 1 - 0.4 = 0.6.$$

Value iteration:

$$V_{k+1}(s) = \max_a \sum_{s'} p(s', r | s, a) [r + \gamma V_k(s')]$$

~~Sabemos que, en este caso, buscar la máxima $p(s'|s)$ no se puede desacoplar: $p(s', r | s, a) = p(s') \cdot p(r | s, a)$?~~



~~$$p(s'=100, r=1 | s=54, a=46) = \dots$$~~

Imaginemos que tenemos q estados: $1, \dots, q$ y
 $a \in \{0, 1, \dots, \min(s, 10-s)\}$

~~$$p(s'=1, r=0 | s=1, a=0) = 1.0$$~~

~~$$p(s'=2, r=0 | s=1, a=1) = 0.4$$~~

~~$$p(s'=0, r=0 | s=1, a=1) = 0.6$$~~

~~$$p(s'=3, r=0 | s=1, a=2) =$$~~

~~$$p(s'=0$$~~

~~$$p(s'=3, r=0 | s=2, a=1) = 0.4$$~~

~~$$p(s'=1, r=0 | s=2, a=1) = 0.6$$~~

~~$$p(s'=4, r=0 | s=2, a=2) = 0.4$$~~

~~$$p(s'=\emptyset, r=0 | s=2, a=2) = 0.6$$~~

según Sutton, la lógica detrás es: dado que la moneda está trucada en contra del agente, lo mejor es apostar lo mínimo posible. cuando llega a un 50% del capital, entonces lo mejor es apostar todo. Pero si viene 51%, puede hacerlo mejor: puede apostar un euro y así ponerse en 52€, con lo que seguiría teniendo la opción de apostar 50 y ganar todo si hubiera necesidad. lo mismo para 53, 54, 55, ... cuando llega a 75, puede apostar 25 y ganar ~~todo~~ en una jugada. Pasados los 75, de nuevo se aplica la lógica anterior: si aposte 1, se pone en 76, con lo que puede ir ganando mientras mantiene los 25€ salvos. para una jugada en un golpe.

~~con respecto a las~~
 Esto explica los cambios bruscos. con respecto a las primas de Sierra, diría que cualquier política con la misma integral bajo esas precisiones de curva funciona?

4.10 - ¿Cuál es el análogo de value iter. update usando $q_{k+1}(s, a)$?

$$q_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_k(s', a')]$$

Ejercicio 4.9: Programming. Implement Value Iteration for gambler's problem...

~~la ecuación~~

El pseudocódigo de Value Iteration es:

Parámetro (γ):

Inicializa $V(s)$, $\forall s \in \mathcal{S}^+$, arbitrariamente excepto $V(\text{terminal}) = 0$.

loop:

$\Delta \leftarrow 0$

loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

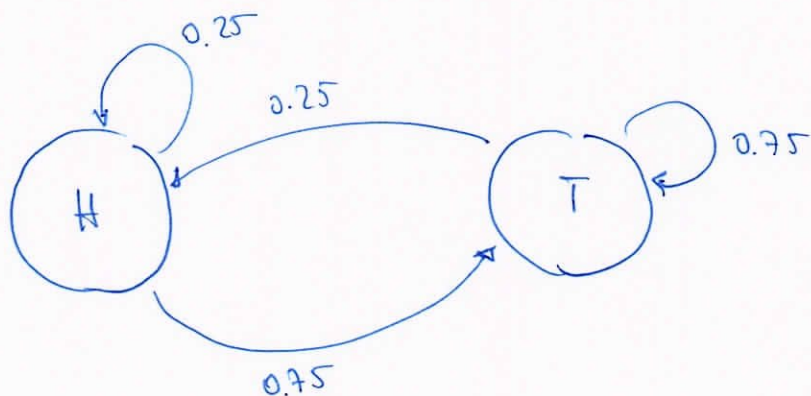
$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

output $\pi \approx \pi^*$:

$\pi(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

En nuestro caso, $\gamma = 1$. También conocemos las probs. de cara y cruz:



$$s' = s + a$$

Si sale cara, ~~QQQ~~. ~~S~~ y $r=0$. Si sale cruz, $s' = s - a$ y $r=0$.

$r=1$ sólo si $s=100$.

Por tanto,

Vamos a tratar de especificar esto en términos de probabilidad $p(s', r | s, a)$:

~~Sabemos que $V(s')$ también queda~~

1) Sabemos que

$$r(s', a, s) = \frac{\sum_r r \cdot p(s', r | s, a)}{p(s' | s, a)}$$

Y por tanto, $\sum_r r \cdot p(s', r | s, a) = r(s', a, s) \cdot p(s' | s, a)$

2) También sabemos que $r(s', a, s)$ en realidad sólo es función de s' : $r(s', a, s) \equiv r(s') = \begin{cases} 0 & \text{if } s' \neq 100 \\ 1 & \text{else} \end{cases}$

3) Por otra parte, sabemos que $p(s' | s, a) = 0.25$ si $s' = s + a$ y $s' = s - a$.
O sea, que dado s , sabemos que $a \in [0, \dots, \min(s, 100-s)]$, y
dados s y a , tenemos:

$$p(s' = s + a | s, a) = 0.25$$

$$p(s' = s - a | s, a) = 0.25$$

Finalmente,

4) También sabemos que $\sum_{s'} p(s', r | s, a) V(s') = \sum_{s'} p(s' | s, a) V(s')$

Por tanto, podemos reescribir la ec. de actualización para $V(s)$ como

$$V(s) \leftarrow \max_a \left[\sum_{s'} r(s', a, s) p(s' | s, a) + \sum_{s'} p(s' | s, a) V(s') \right]$$

o lo que es lo mismo:

$$V(s) \leftarrow \max_a \left[\sum_{s'} p(s'|s,a) (r(s') + V(s')) \right]$$

Dados s y a , la suma anterior es de sólo dos sumandos:

~~$$V(s) \leftarrow \max_a \left[0.25 (r(s+a) + V(s+a)) + 0.75 (r(s-a) + V(s-a)) \right]$$~~

~~$$\max_a \left[0.25 (r(s+a) + V(s+a)) + 0.75 (r(s-a) + V(s-a)) \right]$$~~

Valor calculado para " a ".

Dado s ,

Hay entonces que almacenar estos valores de candidatos a $V(s)$, cada uno asociado a las acciones posibles, y seleccionar el máximo como $V(s)$.

Otra forma de hacer esto es ir guardando, en cada estado, el máximo y ir substituyéndolo. En cada estado, este valor se

renueva:

~~$$\max_a$$~~

~~$$\max_a$$~~
$$V_max_a = -1000$$

$$A = 0.25 \cdot () + 0.75 \cdot ()$$
~~$$\max_a$$~~

~~$$V_max_a = \max(V_max_a, A)$$~~

$$\text{if } A > V_max_a$$

$$V_max_a \leftarrow A$$

~~$$\max_a$$~~

Por tanto, el algoritmo completo queda:

def value-iteration (θ):

$S = \text{~~100~~ [0]} [0, \dots, 100]$

$r = [0, \dots, 1]$

$\vec{V}(\vec{S}) = [0, \dots, 0]$

do

$\Delta \leftarrow 0$

for s_0 in ~~\vec{S}~~ \vec{S}

$v \leftarrow \vec{V}(s_0)$

$v_max_a \leftarrow -1000$

for a in $A(s)$:

$A \leftarrow 0.25 f(s,a) + 0.75 f(s-a)$

if $A > v_max_a$

$v_max_a \leftarrow A$

$\vec{V}(s_0) \leftarrow v_max_a$

$\Delta \leftarrow \max(\Delta, |V - \vec{V}(s_0)|)$

while $\Delta \geq \theta$

$S_0 \equiv S$

$\vec{S} \equiv \text{states}$

$\Delta = \text{~~delta~~ delta}$

$\theta = \text{theta}$

$\Delta(s) = \text{~~available-action~~ get-action-space(s)}$

Bellman update

def $A(s)$:

return $[i \text{ for } i \text{ in } (0, \dots, \min(s, 100-s))]$

Podemos también encapsular la parte del bucle ~~en~~ ~~en~~ ~~en~~ $A(s)$ como:

$\vec{V}(s_0) \leftarrow \text{~~return~~ Bellman-update}(s, a)$, and return v_max_a .