

## Capítulo 5: Ejercicios.

Ej. 5.1: Considera los diagramas de la derecha en la figura 5.1.

Diagrama arriba derecha: el valor ~~del estado~~ de los estados en las últimas dos filas crece con independencia de lo que esté mostrando el dealer.

Cuando el jugador tiene 21 puntos, ha ganado o bien ha empatado, ~~cuando el dealer muestra~~ pero nunca podrá haber perdido.

Si el dealer muestra un As, entonces ~~si está a~~ boca abajo puede tener entre 1 y 10. Si tiene 10, habrá empate. En cualquier otro caso, el jugador gana. Si tiene una carta boca arriba que no es as ni 10, entonces habrá ganado el jugador. Si tiene boca arriba una carta con valor de 10, entonces ~~ganará~~ <sup>habrá empate</sup> si ~~no~~ <sup>si</sup> la boca abajo hay un as. Como hay menos ases que cartas que valen 10, la ~~situación~~ el estado  $(21, A)$  es menos ventajoso ~~que el~~ para el jugador que el estado  $(21, 10)$ .

Si el jugador tiene 20, entonces puede ganar, perder o empatar.

~~En el estado~~  $(20, A)$ , sólo podrá perder si la otra carta es 10, y ~~empatará si~~ <sup>sólo</sup> ~~podrá empatar~~ ~~(no hay x/x+x=20)~~

~~si lo que hay boca abajo es un 9. En el resto de casos,~~  
gana el jugador.

Si tenemos  $(20, 10)$ , entonces el empate sólo se da si hay otro 10 y perderá sólo si hay un as. Este estado es por tanto preferible a  $(20, A)$  para el jugador. Tanto  $(20, 10)$  Además,  $(20, 10) < (21, 10)$  y  $(20, A) < (21, A)$ .



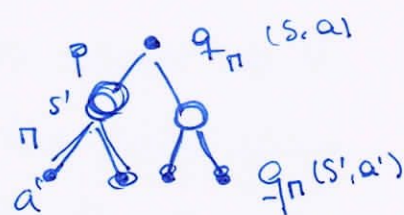
Cuando el jugador tiene en total 19, entonces la política dicta que vuelva a pedir una carta. Esto aumentará drásticamente la probabilidad de pasarle y de que gane el dealer.

5.2. Supón ~~que~~ que utilizamos every-visit MC en vez de first-visit MC en la tarea de blackjack. ¿Cambiará mucho la función? <sup>no estado</sup>

En blackjack, los estados no se repiten. La primera vez que se visita  $r$  es toda las veces que se visita, ~~por~~ por lo que en este caso los resultados serían los mismos.

5.3. ¿Backup diagram for MC estimation of  $q_\pi$ ?

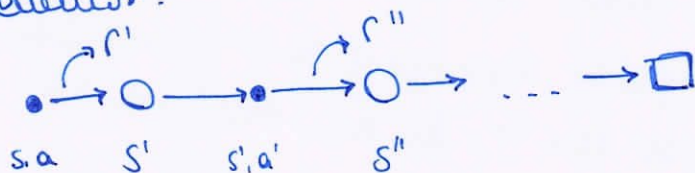
Para calcular  $q_\pi$ , usamos que:



$$\begin{aligned} q_\pi &= \mathbb{E}_\pi [G_t | S_t = s, A_t = a] = \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] = \\ &= \mathbb{E}_\pi [R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}_\pi [\mathbb{E}_\pi [G_{t+1} | S_{t+1} = s', A_{t+1} = a']] = \\ &= \mathbb{E} [R_{t+1} | S_t = s, A_t = a] + \gamma \mathbb{E}_\pi [q_\pi(s', a') | S_t = s] \end{aligned}$$

Por tanto, si generamos un episodio  $S_0, A_0, R_1, A_1, R_2, \dots, R_T$ ,

tenemos:



5.4. : The pseudocode for Monte Carlo ES is inefficient...

En la sección 2.4 habíamos visto que se puede escribir un promedio de forma incremental:

$$Q_n = \frac{\sum_{i=1}^n R_i}{n} = \frac{1}{n} \sum_{i=1}^n R_i = \frac{1}{n} R_1 + \frac{1}{n} [R_1 + \sum_{i=1}^{n-1} R_{i+1}]$$

$$= \frac{1}{n} [R_n + \sum_{i=1}^{n-1} R_i] = \frac{1}{n} R_n + \frac{(n-1)}{n} Q_{n-1} =$$

$$= \frac{1}{n} R_n + \frac{n-1}{n} Q_{n-1} = \frac{1}{n} R_n + \frac{n-1}{n} Q_{n-1} = \frac{Q_{n-1}}{n} =$$

$$= Q_{n-1} + \frac{1}{n} [R_n - Q_{n-1}]$$

$$= Q_{n-1} + \frac{1}{n} [R_n - Q_{n-1}]$$

donde  $R_n \equiv \text{Returns}(S_t, A_t)$  en la  $i$ -ésima iteración.

Por tanto, haríamos, en vez de "append G to returns ...", lo siguiente:

• Mantenemos una cuenta del episodio en que estamos ( $n$ ).

~~Reinicia~~  $n \leftarrow n+1$

~~Q~~  $\leftarrow Q + \frac{1}{n} [G - Q]$

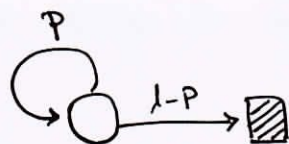
} Esto lo hacemos debajo del "if".

Es decir,  $n \leftarrow n+1$  no se ejecuta

si no es la 1ª vez que  $(S_t, A_t)$  aparece en la secuencia!



Ejercicio 5.5. sea un MDP:



$$r = 1.; \gamma = 1.$$

What are the first-visit and every-visit estimators of the value of the non-terminal state?

Para el caso first-visit, tendríamos  $V(s) = \mathbb{E}[R] = 10$ .

Para el caso every-visit, tenemos  $V(s) = \mathbb{E}[1, 2, 3, \dots, 10] = 5.5$

Ejercicio 5.6. ¿cuál es la ec. análoga a 5.6 para  $Q(s, a)$ , suponiendo los datos los recibiendo siguiendo la política  $b$ ?

Tenemos  $\rho_{t:T-1} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$

$$\mathbb{E}_b[G_t | S_t, A_t] = q_b(s, a). \text{ ¿ } Q_\pi(s, a)?$$

Sabemos que  $q_b(s, a) = \mathbb{E}[R_{t+1} + \gamma \overset{G}{V}^t(s) | S_t = s, A_t = a]$

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

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

$$= \sum_{s', r} p(s', r | s, a) [r + \gamma \rho_{t:T-1} V_b(s')] = \mathbb{E}[r + \gamma \rho_{t:T-1} V_b(s')]$$

$$= \mathbb{E}[R_{t+1} + \gamma \rho_{t:T-1} V_b(S_{t+1}) | S_t = s, A_t = a]$$

Ejercicio 5.6: ¿Etc. análogo a 5.6 para  $Q(s, a)$ ?

En este caso, el retorno  $R_{t+1}$  no viene dado por la estrategia  $b$ , sino por la elección libre de "a". Por tanto, no debemos multiplicar a  $G_t$  por  $\frac{\pi(A_t | S_t)}{b(A_t | S_t)}$ , sino comenzar a usar  $q$  en  $t+1$ :

$$Q(s, a) \doteq \frac{\sum_{t \in \tau(s, a)} \rho_{t+1:T(t)-1} G_t}{\sum_{t \in \tau(s, a)} \rho_{t+1:T(t)-1}}$$

donde  $\tau(s, a)$  ahora contiene los xltos (temporales) en los que se dieron el par estado-acción  $(s, a)$ . (Ref: Bryan Hayder).

Ejercicio 5.7: ¿Por qué creció y luego decreció el error en weighted-importance sampling aplicado a la figura 5.3?

Todas las trayectorias en las que existe alguna acción prohibida en  $\pi$  tienen  $p=0$  y por tanto contribuyen con valor "0" a  $V(s)$ .

Con pocos episodios, ~~se ve~~ ~~se ve~~ bajo el número de trayectorias contribuyendo con ~~el~~ ~~el~~  $p \cdot G \neq 0$  ~~se ve~~ bajo, y  $V(s) \approx 0$ , que es cercano al valor real de  $V(13)$ .

Conforme aumenta el # de episodios, también hay un mayor de trayectorias con  $p \neq 0$  y ~~se ve~~  $V(13)$  aumentará. En este punto aumenta el error ~~se ve~~. Este crecimiento se detiene y reversa conforme la varianza comienza a decrecer.



5.8. los resultados en el ~~ej~~ ejemplo 5.5 mostrados en la figura 5.4 usan first-visit MC. Supón que usaramos every-visit MC. ¿Seguiría siendo la variante infinita?

Tanto en este caso como en first-visit, nuestro objetivo es calcular el valor esperado del retorno escalado y el cuadrado.

~~Introducción~~

El valor esperado de una variable aleatoria  $x$  puede obtenerse, o bien con la media muestral (es decir, generando muchos episodios, calculando  $G_t$  en cada episodio y ~~dividiendo~~ luego sumando los  $G_t$ 's obtenidos y dividiendo entre # de ~~ep~~ muestras), o bien ~~con~~ de forma exacta. Para obtenerla de forma exacta, es necesario multiplicar la probabilidad de cada posible valor de la variable aleatoria por dicho valor, y después sumar.

Dado que conocemos la probabilidad de cada valor porque conocemos la dinámica del MDP, calculamos  $\mathbb{E}_b[X^2]$  como

$$\sum_i p_b(x_i) \cdot x_i^2, \text{ donde } x_i^2 \text{ es la variable aleatoria y } p_b(x_i)$$

es la probabilidad de obtener  $x_i$  si se sigue la dinámica del MDP junto con la estrategia  $b$ :

$$\mathbb{E}_b[X^2] = \sum_i x_i^2 p_b(x_i)$$

$$p_b(x_i) = P_b(\text{trayectoria}) \neq$$

Dada una trayectoria de longitud  $T$ , su probabilidad  $\pi$  a los estados de la trayectoria obteniéndose multiplicando probs. de transición  $p(s'|s, a)$  y prob. de elección de las acciones ejecutadas en la misma:

$$P_b\{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim b\} = \prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)$$



En el caso de  $\text{every-time visit}$ , si estamos en una trayectoria de longitud  $T$ , necesariamente habremos pasado  $T-1$  veces por " $s$ " (el único estado del problema). Por lo tanto,  $T(s) = \{1, 2, \dots, T-1\}$  y

tenemos

$$V(s) = \frac{\sum_{t=1}^{T-1} p_{t:T(t)-1} G_t}{T-1} \quad T \geq 2$$

Donde  $G_t = G_0 = 1 \quad \forall t$  y  $T(t) = T \quad \forall t$ , con lo que  $V(s)$  queda:

$$V(s) = \frac{\sum_{t=1}^{T-1} p_{t:T-1} G_0}{T-1} = \sum_{t=1}^{T-1} \prod_{k=t}^{T-1} \frac{p(A_k | S_k)}{b(A_k | S_k)} G_0$$

Para obtener  $\mathbb{E}_b[V(s)^2]$ , habrá que multiplicar estos valores por la probabilidad de cada trayectoria. Dada una trayectoria de longitud  $T$ , decíamos que

$$\Pr\{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim b\} = \prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

Por lo tanto, tendremos que  $\mathbb{E}_b[V(s)]^2 =$

$$= \sum_{T=2}^{\infty} \underbrace{\prod_{k=1}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)}_{\text{prob. trayectoria de longitud } T} \cdot \frac{1}{(T-1)^2} \left( \sum_{t=1}^{T-1} p_{t:T-1} \right)^2 (*)$$

$t=1 \quad t=2$   
 $\square \rightarrow \square$   
 prob. trayectoria de longitud  $T$

Ej.:  $T=2 \quad (k=1)$

$$b(A_1 | S_1) p(S_2 | S_1, A_1) \cdot \left(\frac{1}{1}\right)^2 \cdot e^2 = \frac{0.1}{2} \cdot \left(\frac{1}{0.5}\right)^2$$

$T=3 \quad (k=2)$

$$b(A_1 | S_1) p(S_2 | S_1, A_1) \cdot b(A_2 | S_2) p(S_3 | S_2, A_2) \cdot \left(\frac{1}{2}\right)^2 \cdot \left[\left(\frac{1}{0.5}\right)^2 + \left(\frac{1}{0.5}\right)^2\right] = \dots$$



$$\frac{1}{0.5^2} = 2^2 = 4$$

$$\dots \left(\frac{1}{2}\right)^2 \cdot 0.1 \cdot 0.9 \cdot \left(\frac{1}{2}\right)^2 \left[2 \cdot \frac{1}{0.5^2}\right]$$

Este resultado Nótese que  $P_{t:T-1} = \prod_t^{T-1} 2 = 2^{(T-t)}$

Manipulando la expresión anterior, queda:

$$\sum_{T=2}^{\infty} \prod_{k=1}^{T-1} \frac{1}{2} P(S_{k+1}|S_k, A_k) \left(\frac{1}{T-1}\right)^2 \left(\sum_{t=1}^{T-1} 2\right)^2 \approx \frac{2}{3}$$

... cuyo límite en el infinito, es también  $\frac{2}{3}$  infinito  
 $= \frac{2}{3}$  por el término exponencial?

Ejercicio 5.11. En el algoritmo, cabría esperar ver el ratio

$$\frac{\pi(A_t|S_t)}{b(A_t|S_t)}, \text{ pero en su lugar aparece } \frac{1}{b(A_t|S_t)} \text{ ¿Por qué?}$$

~~Estados considerando todos~~

Porque hacemos a  $\pi(s)$  una política determinista. Concretamente,

$$\pi(\arg \max_a Q(s_t, a) | S_t) = 1.0 \text{ y si } A_t \neq \pi(S_t), \text{ entonces}$$

$\pi(A_t|S_t) = 0$ . Por tanto, para las acciones que maximizan  $Q(s_t, a)$  (que son a las que se llega cuando con los que se llega a la línea de código que actualiza a  $w$ ),  $\pi(a|S_t) = 1$ , de ahí la actualización

$$W \leftarrow W \frac{1}{b(A_t|S_t)}$$

Si  $x$  elige una acción  $A_t$  que es imposible en  $\pi$ , los cálculos se frenan y el episodio no avanza (aunque sí sirve para calcular  $Q$  y actualizar  $Q(s_t, A_t)$ )



Eq. 5.9: Modifica el algoritmo de first-visit MC para evaluación de estrategias usando la implementación incremental de la sección 2.4.

Offline:

Input: an arbitrary target policy  $\pi$

### Inhibitor:

$$V(s) \in \mathbb{R} \text{ arbitrarily, for all } s \in \mathcal{S}.$$
$$h(s) = 0 \quad \forall s \in S.$$

*[Signature]*

loop forever (for each episode):

Generalize episode following  $b: S_0, A_0, R_1, S_1, A_1, R_2, S_2, \dots, S_{T-1}, A_{T-1}, R_T$

$$G \rightarrow O$$
$$w \leftarrow 1$$

Loop for each step of episode:  $t = T-1, \dots, 0$ :

$$\begin{aligned} G &\leftarrow G + \gamma R_{t+1} \\ W &\leftarrow \frac{\pi(A_t | S_t)}{b(A_t | S_t)} W \end{aligned}$$

Unless  $S_t$  appears in  $S_0, \dots, S_{t-1}$ :

$$u(S_t) \leftarrow u(S_t) + 1$$
$$V(S_t) \leftarrow V(S_t) + \frac{W}{h(S_t)} [G - V(S_t)]$$

online :

Input:  $\pi$  (a predictor).

Inhibitor:

realize:  
 $v(s) \in \mathbb{R}$  arbitrarily for all  $s \in S$

$$M(S) = 0 \quad \forall S \in \mathcal{S}$$

loop break :

Generalized episode

for each  $t$  in  $\mathbb{Z}$  <sup>step in episode:</sup>  $t = T-1, \dots, 0$ :

$$G \leftarrow G + \gamma R_{t+1}$$

unless  $S_t$  in  $S_0, \dots, S_{t-1}$ :

~~68~~ - 69 -



$$n(S_t) \leftarrow n(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) - \frac{1}{n(S_t)} [G - V(S_t)]$$

Ej. 5.10: Deriva la regla 5.8 de 5.7.

Sabemos del capítulo 2 que dado un valor  $G = \frac{\sum R_i}{n}$

$$E[G_t | S_t] = E\left[\underbrace{\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}}_{G_t} | S_t\right] \approx \sum_{i=1}^N \frac{G_{t,i}}{N}$$

donde  $G_i \equiv$  cada uno de los retornos obtenidos al empezar en  $S_t$  en la muestra y  $N \equiv$  n° veces en que se ha observado  $S_t$  en la muestra.

$$\text{Por tanto, tenemos: } V_{n+1} = \frac{1}{n} \sum_{i=1}^n G_i = \frac{1}{n} \left( G_n + \sum_{i=1}^{n-1} G_i \right) =$$

$$= \frac{1}{n} \left( G_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} G_i \right) = \frac{1}{n} \left( G_n + (n-1) V_n \right) =$$

$$= \frac{1}{n} (G_n + nV_n - V_n) = V_n + \frac{1}{n} [G_n - V_n]$$

$$\text{En 5.7, lo que tenemos es } V_{n+1} = \frac{\sum_{k=1}^n W_k G_k}{\sum_{k=1}^n W_k}$$

Sabemos por tanto que

$$V_n = \frac{\sum_{k=1}^{n-1} W_k G_k}{\sum_{k=1}^{n-1} W_k} \rightarrow \sum_{k=1}^{n-1} W_k G_k = \sum_{k=1}^{n-1} W_k V_n (*)$$



Además,  $\sum_{k=1}^n W_k G_k = \sum_{k=1}^{n-1} W_k G_k + W_n G_n$

Sustituyendo (\*) en esta última igualdad:

$$\sum_{k=1}^n W_k G_k = \sum_{k=1}^{n-1} W_k V_n + W_n G_n \quad (a)$$

También sabemos que  $\sum_{k=1}^{n-1} W_k = \sum_{k=1}^n W_k - W_n \quad (**)$

Sustituyendo (\*\*) en (a):

$$\sum_{k=1}^n W_k G_k = V_n \cdot \left[ \sum_{k=1}^n W_k - W_n \right] + W_n G_n \quad (b)$$

Manipulando:

$$\sum_{k=1}^n W_k G_k = \sum_{k=1}^n W_k V_n + W_n [G_n - V_n] \quad (c)$$

Y ahora sustituimos (c) en el ~~deno~~ numerador de 5.7:

$$V_{n+1} = \frac{\sum_{k=1}^n W_k G_k}{\sum_{k=1}^n W_k} = \frac{\sum_{k=1}^n W_k V_n + W_n [G_n - V_n]}{\sum_{k=1}^n W_k} =$$

$$= V_n + \frac{W_n}{\sum_{k=1}^n W_k} [G_n - V_n] = V_n + \frac{W_n}{C_n} [G_n - V_n] \quad (5.8)$$

donde  $C_n = \sum_{k=1}^n W_k$