

Ejercicios del Capítulo 2.

Ej. 2.1.

Los dos métodos ϵ -greedy tenderán a encontrar ~~el~~ ^{el valor} real de $q(a)$:

$$\lim_{t \rightarrow \infty} Q(a) = q(a)$$

~~Sin embargo,~~ ϵ .

Por tanto, el ~~comportamiento~~ comportamiento a largo plazo de ambos será el mismo (no siendo éste el caso del algoritmo sin exploración, que puede quedar atascado en un valor de Q y no seguir mejorando).

El método con $\epsilon = 0.1$ usará a la larga la acción óptima en 90% de las veces, mientras que el de $\epsilon = 0.01$ lo hará en 99% de las veces, lo que supone una mejora del 10%.

Ej. 2.2.

Habíamos llegado a una expresión como

$$Q_{n+1} = \alpha_n R_n + \alpha_{n-1} (1 - \alpha_n) R_{n-1} + (1 - \alpha_n)(1 - \alpha_{n-1}) Q_{n-1}$$

Podemos extenderlo fácilmente un paso más:

$$Q_{n+1} = \alpha_n R_n + \alpha_{n-1} (1 - \alpha_n) R_{n-1} + \alpha_{n-2} (1 - \alpha_{n-1})(1 - \alpha_{n-2}) R_{n-2} + \\ + (1 - \alpha_n)(1 - \alpha_{n-1})(1 - \alpha_{n-2}) Q_{n-2}$$

El término en R_i queda $\sum_i \alpha_i R_i \prod_{k=i+1}^n (1-\alpha_k)$

Y el término en Q :

$$\prod_{i=1}^n (1-\alpha_i) Q_1$$

Y por tanto:

$$Q_{n+1} = \prod_{i=1}^n (1-\alpha_i) Q_1 + \sum_i \alpha_i R_i \prod_{k=i+1}^n (1-\alpha_k)$$

Ej. 2.5

~~Después la primera parte,~~

$$Q_{k+1} = Q_k + \alpha [R_k - Q_k]$$

Si partimos de Q_1 mucho más alto que los demás la media de recompensas que se va a obtener, tardaremos en converger a un valor de Q_k :

$$Q_2 = Q_1 + \alpha [R_1 - Q_1]$$

$$A_{t+2} = \arg \max_a Q_t(a)$$

El valor seleccionado de Q_1 y, concretamente, lo cercano a \bar{R} que sea, determinará si este método converge rápido o lento.

Supongamos Q_1 fijo, con $Q_1 = 5$.

~~Si α es pequeño~~

Supongamos que calculamos Q_2 de $\{a_1, \dots, a_n\}$ y encontramos la acción óptima.

En ese punto, el algoritmo actuará de forma "greedy" y aspirará seleccionándola, convergiendo a un valor real (probablemente menor que 5). En ese punto, el resto de acciones tendrán aún un valor estimado superior a 5, y ~~entonces~~ debido a falta de conocimiento del algoritmo, por lo que el agente acabará de acción a una subóptima.

Ej. 2.6.

Caso A: $\begin{cases} a_1: 0.1 (50\%) & \text{y } 0.2 (50\%) \\ a_2: 0.2 (50\%) \end{cases}$

Caso B: $\begin{cases} a_1: 0.9 (50\%) \\ a_2: 0.8 (50\%) \end{cases}$

No sé ~~cuál~~ en qué caso estoy. $\mathbb{E}[R]$? ¿Cómo debo comportarme?

Si estoy en el caso A, lo mejor que puedo hacer es escoger a_2 .

" " " " " B, " " " " " a_1 .

Si juego 100 veces seleccionando " a_1 ", tendré ~~0.5~~ $0.05 + 0.45 = 0.5$

" " " " " a_2 ", " $0.01 + 0.4 = 0.5$

¿Y si selecciono al ~~ale~~ azar?

$$0.5 \cdot 0.5 + 0.5 \cdot [0.5 \cdot 0.1 + 0.5 \cdot 0.2] + 0.5 [0.5 \cdot 0.9 + 0.5 \cdot 0.8] =$$

$$= \frac{1}{2} [0.15] + \frac{1}{2} [0.45 + 0.4] = \frac{1}{2} 0.15 + \frac{1}{2} 0.85 =$$

$$= 0.075 + 0.425 = 0.5$$

Por tanto, de igual lo que haga (siempre a_1 , siempre a_2 o $\frac{1}{2}a_1 + \frac{1}{2}a_2$), siempre obtengo $\mathbb{E}[R]$ de 0.5.

Imaginemos ahora que sabemos si estamos en A o B.

Si estoy en A, debo escoger 2. Si estoy en B, debo escoger 1.

Espero tener $0.5 \cdot 0.2 + 0.5 \cdot 0.9 = 0.1 + 0.45 = 0.55$

Ej. 2.2.

* n-bandit

$$\alpha = \frac{1}{k} \cdot \frac{1}{n}$$

Greedy.

bandit(a) = R

Usar arrays y variables, modo de subscripts.

Recordemos:

$$Q_{u+1} = Q_u + \alpha (R_u - Q_u)$$

$$A_u = \arg \max_a Q_u(a)$$

Variables que guardo: Q, n, A

~~aux~~ ← -1000

$Q \leftarrow 0$

$n \leftarrow 1$
for $u \in [1, 2000]$:

for $a \in A$:

if $Q + \alpha (\text{bandit}(a) - Q) > \text{aux}$

~~aux~~ ← $Q + \alpha (\text{bandit}(a) - Q)$

$Q \leftarrow \text{aux}$

~~aux~~ = -1000

Ej. 2.4: Diseña un ej. para demostrar las dificultades a las que se enfrentan los ~~estos~~ métodos average-sample en problemas no estacionarios.

Usar 10-armed bandit, con $g(a)$ constantes igual y después siendo independientes (random walks).

Usar $\alpha = \frac{1}{k}$ y otro con $\alpha = 0.1$

Usar $\epsilon = 0.1$

Necesito:
 → Una función que genere los valores $g(a_1), \dots, g(a_n)$ en cada instante
 → Una función que genere el ruido.
 → $\text{bandit}(a)$ (combina ambos).
 → Una función que escoja la acción
 → greedy
 → pero con ϵ -exploración

El pseudocódigo es
 → $Q_{n+1} = Q_n + \alpha (R_n - Q_n)$ ($\alpha = \frac{1}{k}$)
 (Es decir, como en el problema 2.2).

He hecho un programa que se comporta de manera similar a las curvas del libro. No obstante, a continuación transcribo y desmenuzo el esquema pedido por NeuroMatch Academy (curra carpete doctorado > RL):

1° Définir une fonction qui sélectionne la action de max value
con prob. $1-\epsilon$, et une action aléatoire con prob. ϵ .

def epsilon-greedy (q, epsilon):

Args:

q (ndarray): array de valeurs de la s u-actiones.

epsilon (float): prob. de ~~select~~ explorer.

si random > epsilon:

selection: $a = \operatorname{argmax}(q)$

else:

$a = \text{selection aléatoire} \Rightarrow a = \text{np.random.choice}(len(q))$

return a

def update_action_value (q, action, reward, alpha):

Args:

q (ndarray).

action (int)

reward (float)

alpha (float)

Returns:

float: the updated value for the selected action.

$$\text{value} = q[\text{action}] + \alpha \cdot (r - q[\text{action}])$$

return value.

Finalmente, ~~la~~ el programa `ppal`:

`def multiarmed_baudit(n_arms, epsilon, alpha, n_steps):`

Args:

`n_arms` (int)

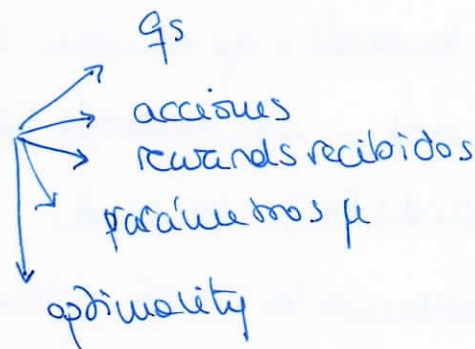
`epsilon` (float)

`alpha` (float)

`n_steps` (int)

Returns:

`dict`: un diccionario con



#1: las acciones tienen valores reales $q(a)$ provenientes de una

$N(0,1)$:

$\mu = \text{np.random.normal}(\text{size} = \text{n_arms}) \rightarrow \text{vector } \vec{\mu}(a)$

los valores calculados $\vec{Q}_t(a)$ se inicializan a ϕ .

Guarden los Q_t de cada iteración para analizar.

También guarden los rewards observados (uno por ~~tiempo~~ iteración),

las acciones seleccionadas (una por iteración) y el # de veces que se escoge una óptima.

inicialmente, todo es ϕ :

$q = \text{np.zeros}(n_arms)$

$q_t = \text{np.zeros}(n_steps, n_arms)$

$\text{rewards} = \text{np.zeros}(n_steps)$

$\text{actions} = \text{np.zeros}(n_steps)$

$\text{optimal} = \text{np.zeros}(n_steps)$

Run the bandit

for t in range(n_steps):

$\text{action} = \text{epsilon-greedy}(q, \text{epsilon})$

$\text{actions}[t] = \text{action}$

$\text{all_rewards} = \text{np.random.normal}(\mu_0)$ ← premios observados. Tenemos media μ y varianza 1

$\text{reward} = \text{all_rewards}[\text{action}]$

$\text{rewards}[t] = \text{reward}$

buscaremos la acción óptima y comprobaremos si coincide

con action

$\text{optimal_action} = \text{np.argmax}(\text{all_rewards})$

$\text{optimal}[t] = \text{action} \text{ } \neq \text{ } \text{optimal_action}$

Update the action value

$q[\text{action}] = \text{update_action_value}(q, \text{action}, \text{reward}, \alpha)$

$q_t[t] = q$

results = { 'q_t' : q_t ,
 'actions' : actions ,
 'rewards' : rewards ,
 'optimal' : optimal ,
 '-' : 16-