

1 → Comment

Initialize

$$\pi(s) \in A(s) \text{ (arbitrarily)}, \text{ for all } s \in S$$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in S, a \in A(s)$

Return $(\text{key}) \rightarrow$ empty list \forall for all $\text{key} \in \text{keys}, \text{key} \in \text{keys}$

① $\text{Count}(s, a) \in W$, initially 0 for all $s \in S, a \in A(s)$
 // maintain count of occurrence of each ^{state}-action pair

Loop forever (for each episode):

Choose $S_0 \in S$, $A_0 \in A(S_0)$ randomly such that all pairs have $p > 0$

Generate episode from S_0, A_0 following $\Pi: S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$$G_1 \leftarrow 0$$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

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

Unless pair S_k, A_k appears in $S_0, A_0, S_1, A_1, \dots, S_{k-1}, A_{k-1}$:

③ ~~if~~ $\text{Count}(S_k, A_k) + 1$ // increase count

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{C(s_t, a_t)} (G_t - Q(s_t, a_t))$$

~~CSP, ACTA~~ ~~Kinase~~ ~~not~~ ~~for~~ ~~this~~ ~~action~~

$$\Pi(s_t) \leftarrow \arg \max_a Q(s_t, a)$$

Explanation

Changes are lines ^{marked} ①, ② and ③

① was originally a list Returns (S_t, A_t) for all states action pairs

② originally calculated the new average of all return for S_t, A_t

$$\text{New Avg} = (\text{Old avg} \times (\text{count}-1) + \text{new value}) / (\text{count})$$
$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{1}{\text{Count}(s_t, A_t)} [r_t - Q(s_t, A_t)]$$