**Off-policy MC control, for estimating $\pi \approx \pi_*$**

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:
$\quad Q(s, a) \in \mathbb{R}$ (arbitrarily)
$\quad C(s, a) \leftarrow 0$
$\quad \pi(s) \leftarrow \operatorname{argmax}_a Q(s, a)$    (with ties broken consistently)

Loop forever (for each episode):
$\quad b \leftarrow$ any soft policy
$\quad$ Generate an episode using $b$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
$\quad G \leftarrow 0$
$\quad W \leftarrow 1$
$\quad$ Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
$\quad\quad G \leftarrow \gamma G + R_{t+1}$
$\quad\quad C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
$\quad\quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} \left[ G - Q(S_t, A_t) \right]$
$\quad\quad \pi(S_t) \leftarrow \operatorname{argmax}_a Q(S_t, a)$    (with ties broken consistently)
$\quad\quad$ If $A_t \neq \pi(S_t)$ then exit inner Loop (proceed to next episode)
$\quad\quad W \leftarrow W \frac{1}{b(A_t | S_t)}$