1 a) weighted Average

$$v_n = \frac{\sum_{k=1}^{n-1} \omega_k G_k}{\sum_{k=1}^{n-1} \omega_k}$$

$$v_{n+1} = \frac{\sum_{k=0}^{n} \omega_k G_k}{\sum_{k=0}^{n} \omega_k}$$

$$= \frac{1}{\sum_{k=0}^{n} \omega_k} \times \left( \text{set } \frac{1}{C_n} \right)$$

$$= \frac{1}{C_n} \left( \omega_n G_n + \sum_{k=0}^{n-1} \omega_k G_k \right)$$

where $\sum_{k=0}^{n-1} \omega_k G_k = \left( \sum_{k=0}^{n-1} \omega_k \right) v_n$

$$= \sum_{k=0}^{n} \omega_k - \omega_n$$

$$\therefore \quad v_{n+1} = v_n + \frac{\omega_n}{\sum_{k=0}^{n} \omega_k} (G_n - v_n)$$

$\boxed{\text{C}_n \text{ can covered in the same method}}$

$$\Rightarrow C_n = \sum_{k=0}^{n} \omega_k = \omega_n - \sum_{k=0}^{n-1} \omega_k$$

$$\therefore \quad v_{n+1} = v_n + \frac{\omega_n}{C_n} [G_n - v_n], \quad n \geq 1$$

1 b) Initialize, for all $s \in S$, $a \in A(S)$

$\quad Q(s, a) \in R$ (arbitrarily)

$\quad C(s, a) \leftarrow 0$

$\quad \pi(s) \leftarrow \text{argmax}_a Q(s, a)$ (with ties broken consistently)

$\quad$ Loop forever (for each episode):

$\quad\quad b \leftarrow$ any soft policy

$\quad\quad$ Generate an episode using $b$: $S_0, A_0, R_1 \dots$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \dots S_{T-1}, A_{T-1}, R_T$

$\quad\quad G \leftarrow 0$

$\quad\quad W \leftarrow 1$

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

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

$\omega = 0 \atop \text{after } A_t \neq \pi(S_t)} \Big\{ \begin{array}{l} C(S_t, A_t) \leftarrow C(S_t, A_t) + W \\ Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \dfrac{W}{C(S_t, A_t)} \big[ G - Q(S_t, A_t) \big] \end{array}$

$\quad\quad\quad \pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$ (with ties broken consistently)

$\quad\quad\quad$ if $A_t \neq \pi(S_t)$ then exit inner loop

$\omega = 0 \atop \text{for all } t \atop \text{after } A_t \neq \pi(S_t)} \Big\{ W \leftarrow W \dfrac{1}{b(A_t | S_t)}$

---

$\quad$ The use of $1/b(A_t | S_t)$ instead of $\dfrac{\pi(A_t | S_t)}{b(A_t | S_t)}$ in the update of $W$ is correct because $\left[ \dfrac{\pi(A_t | S_t)}{b(A_t | S_t)} \right] W$ is itself an accumulation of the importance-sampling ratios from the beginning of the episode to the time $t$.

and adjusts the update at each time step. The policy $\pi$ comes into play by determining which actions contributed to the update, if $\pi(A_t | S_t)$ is zero, then the update would be zero, effectively filtering out the updates from actions not supported by the target policy $\pi$, even though the data is generated by following behaviour policy $b$.

3a) The agent must have moved to left terminal from state A receiving no reward, and since the terminal state has a value of 0:

$$V(A) \leftarrow 0.5 + 0.1 [0 + 1(0) - 0.5]$$

$$= 0.5 - 0.05 = 0.45$$

All the other states were updated as well having the same reward and same initialization the error update in TD was very small

Ex :-

$$v(B) \sim 0.5$$

$$V(c) \leftarrow 0.5 + 0.1 [0 + 1 \cdot 0.5 - 0.5]$$

$$= 0.5$$