## Expectation

Input: Behaviour policy $\pi$, target policy $\pi_i$, discount factor gamma, initial state $S_o$.

$R \leftarrow [\ ]$, discounted rewards $\leftarrow 0$, ~~$\alpha \leftarrow [A \leftarrow [\ ] A \leftarrow [\ ]$~~ $Q \leftarrow [\ ][\ ]$

for $j = 0, 1, 2 \ldots t$ do: // $t$ is terminal state

    Take an action $a_j$ at $S_j$ according to $\pi$, collect reward $r_j$ and move on to $S_{j+1}$.

    Store $R[j] \leftarrow r_j$

    Store $A[j] \leftarrow a$

    Store $Q[j][a] = \pi(a | s_j)$

$s = 0$

for $r$ in ~~R do~~ reverse(R) do:

    if $r$ is last element:

        $s = 0$

    else:

        $s = r + \gamma * s$

    ~~Ap~~ //Create a new list R1,

    R1.append(s) at $0^{th}$ position.

~~for $j$ in A do:~~

~~$L = \frac{\pi_i}{\pi}$~~

sum = 0

for $j$ in A do:

    $L_\pi = \dfrac{\pi_i(a|s_j)}{Q[j][A[j]]} \times R_1[j]$

    sum += $L_\pi$

$$surr = sum - \frac{(4 \times max(R_1) \times gamma) \times d\_kl(\pi, \pi_i)}{(1 - \gamma)^2}$$

return surr