# Implementation: Sarsa(0)

The pseudocode for Sarsa (or Sarsa(0)) can be found below.

<div style="border: 2px solid #00bcd4;">

## TD Control: Sarsa(0)

**Input:** policy $\pi$, positive integer $num\_episodes$, small positive fraction $\alpha$
**Output:** value function $Q$ ($\approx q_\pi$ if $num\_episodes$ is large enough)
Initialize $Q$ arbitrarily (e.g., $Q(s,a) = 0$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$, and $Q(terminal\text{-}state, \cdot) = 0$)
**for** $i \leftarrow 1$ **to** $num\_episodes$ **do**
$\quad \epsilon \leftarrow \frac{1}{i}$
$\quad$ Observe $S_0$
$\quad$ Choose action $A_0$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
$\quad t \leftarrow 0$
$\quad$ **repeat**
$\quad\quad$ Take action $A_t$ and observe $R_{t+1}, S_{t+1}$
$\quad\quad$ Choose action $A_{t+1}$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
$\quad\quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$
$\quad\quad t \leftarrow t + 1$
$\quad$ **until** $S_t$ is terminal;
**end**
**return** $Q$

</div>

Sarsa(0) is **guaranteed to converge** to the optimal action-value function, as long as the step-size parameter $\alpha$ is sufficiently small, and the **Greedy in the Limit with Infinite Exploration (GLIE)** conditions are met. The GLIE conditions were introduced in the previous lesson, when we learned about MC control. Although there are many ways to satisfy the GLIE conditions, one method involves gradually decaying the value of $\epsilon$ when constructing $\epsilon$-greedy policies.

In particular, let $\epsilon_i$ correspond to the $i$-th time step. Then, if we set $\epsilon_i$ such that:

- $\epsilon_i > 0$ for all time steps $i$, and
- $\epsilon_i$ decays to zero in the limit as the time step $i$ approaches infinity (that is, $\lim_{i \to \infty} \epsilon_i = 0$),

then the algorithm is guaranteed to yield a good estimate for $q_*$, as long as we run the algorithm for long enough. A corresponding optimal policy $\pi_*$ can then be quickly obtained by setting $\pi_*(s) = \arg\max_{a \in \mathcal{A}(s)} q_*(s, a)$ for all $s \in \mathcal{S}$.

Implementation

If you'd like to reference the pseudocode while working on the notebook, you are encouraged to open **this sheet** in a new window.

Feel free to check your solution by looking at the corresponding section in

NEXT

If you'd like to reference the pseudocode while working on the notebook, you are encouraged to open **this sheet** in a new window.