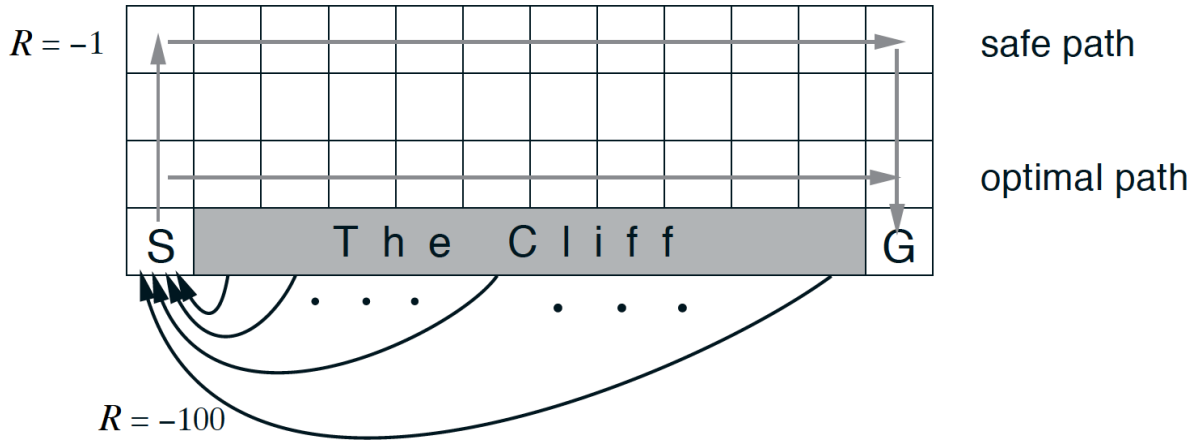# 1.9 Temporal Difference Methods Summary



**The cliff-walking task (Sutton and Barto, 2017)**

## Temporal-Difference Methods

- Whereas Monte Carlo (MC) prediction methods must wait until the end of an episode to update the value function estimate, temporal-difference (TD) methods update the value function after every time step.

## TD Control

- **Sarsa(0)** (or **Sarsa**) is an on-policy TD control method. It is guaranteed to converge to the optimal action-value function $q_*$, as long as the step-size parameter $\alpha$ is sufficiently small and $\epsilon$ is chosen to satisfy the **Greedy in the Limit with Infinite Exploration (GLIE)** conditions.

---

**Algorithm 13:** Sarsa

---

**Input:** policy $\pi$, positive integer $num\_episodes$, small positive fraction $\alpha$, GLIE $\{\epsilon_i\}$
**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**
    $\epsilon \leftarrow \epsilon_i$
    Observe $S_0$
    Choose action $A_0$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
    $t \leftarrow 0$
    **repeat**
        Take action $A_t$ and observe $R_{t+1}, S_{t+1}$
        Choose action $A_{t+1}$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
        $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$
        $t \leftarrow t + 1$
    **until** $S_t$ is terminal;
**end**
**return** $Q$

---

- **Sarsamax** (or **Q-Learning**) is an off-policy TD control method. It is guaranteed to converge to the optimal action value function $q_*$, under the same conditions that guarantee convergence of the Sarsa control algorithm.

---

**Algorithm 14:** Sarsamax (Q-Learning)

---

**Input:** policy $\pi$, positive integer $num\_episodes$, small positive fraction $\alpha$, GLIE $\{\epsilon_i\}$
**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**
    $\epsilon \leftarrow \epsilon_i$
    Observe $S_0$
    $t \leftarrow 0$
    **repeat**
        Choose action $A_t$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
        Take action $A_t$ and observe $R_{t+1}, S_{t+1}$
        $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$
        $t \leftarrow t + 1$
    **until** $S_t$ *is terminal*;
**end**
**return** $Q$

---

- **Expected Sarsa** is an on-policy TD control method. It is guaranteed to converge to the optimal action value function $q_*$, under the same conditions that guarantee convergence of Sarsa and Sarsamax.

---

**Algorithm 15:** Expected Sarsa

---

**Input:** policy $\pi$, positive integer $num\_episodes$, small positive fraction $\alpha$, GLIE $\{\epsilon_i\}$
**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**
    $\epsilon \leftarrow \epsilon_i$
    Observe $S_0$
    $t \leftarrow 0$
    **repeat**
        Choose action $A_t$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
        Take action $A_t$ and observe $R_{t+1}, S_{t+1}$
        $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a) - Q(S_t, A_t))$
        $t \leftarrow t + 1$
    **until** $S_t$ *is terminal*;
**end**
**return** $Q$

---

## Analyzing Performance

---

- On-policy TD control methods (like Expected Sarsa and Sarsa) have better online performance than off-policy TD control methods (like Q-learning).
- Expected Sarsa generally achieves better performance than Sarsa.