

Monte Carlo vs Temporal Difference

(an episode can be represented as a character dying and the steps he took)

- with Monte Carlo, we update the value function from a complete episode, and so we use the actual accurate discounted (current) return of this episode

- with TD, we update the value function from a step, so we replace G_t with an estimated return

$$MC: V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

α is learning rate

$$TD: V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Q-Learning

Q-Learning is the algorithm used to train a Q-Function which is an action-value function.

Internally, our Q-Function has a Q-table, a table where each cell corresponds to a state-action pair value.

Example:

S \ A	←	→	↑	↓
S ₀	0	0.1	10	-4
S ₁	0	0	10	0
S ₂	5	-3	0	0.5
S ₃	3	-1	1.5	0
S ₄	1	2	-3	1
S ₅	0	-2	3	0