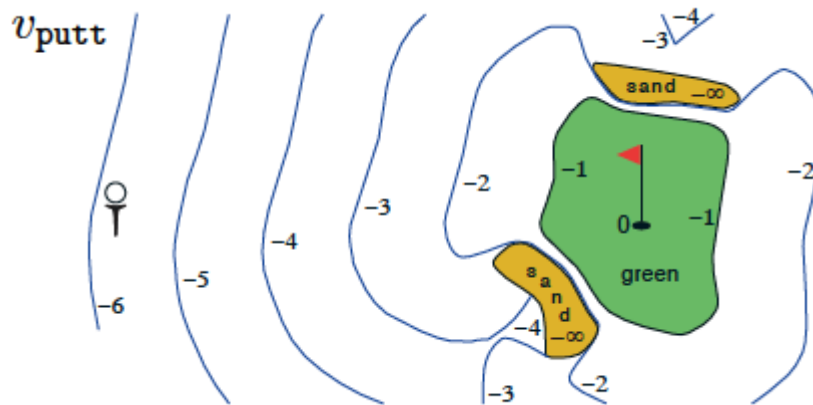# Summary



State-value function for golf-playing agent (Sutton and Barto, 2017)

## Policies

---

- A **deterministic policy** is a mapping $\pi:S\to A$. For each state $s\in S$, it yields the action $a\in A$ that the agent will choose while in state $ss$.
- A **stochastic policy** is a mapping $\pi:S\times A\to[0,1]$. For each state $s\in S$ and action $a\in A$, it yields the probability $\pi(a|s)$ that the agent chooses action $aa$ while in state $ss$.

## State-Value Functions

---

- The **state-value function** for a policy $\pi$ is denoted $v_\pi$. For each state $s\in S$, it yields the expected return if the agent starts in state $ss$ and then uses the policy to choose its actions for all time steps. That is,
$$v_\pi(s)\doteq E_\pi[G_t|S_t=s]$$
We refer to $v_\pi(s)$ as the **value of state** $ss$ **under policy** \pi$\pi$.
- The notation $E_\pi[\cdot]$ is borrowed from the suggested textbook, where $E_\pi[\cdot]$ is defined as the expected value of a random variable, given that the agent follows policy $\pi$.

## Bellman Equations

- The **Bellman expectation equation for** $v_\pi$ is: $v_\pi(s)=E_\pi[R_{t+1}+\gamma v_\pi(S_{t+1})|S_t=s]$.

## Optimality

- A policy $\pi'$ is defined to be better than or equal to a policy $\pi$ if and only if $v_{\pi'}(s)\geq v_\pi(s)$ for all $s\in S$.

- An **optimal policy** $\pi_*$ satisfies $\pi_* \geq \pi$ for all policies $\pi$. An optimal policy is guaranteed to exist but may not be unique.
- All optimal policies have the same state-value function $v_*$, called the **optimal state-value function**.

## Action-Value Functions

- The **action-value function** for a policy $\pi$ is denoted $q_\pi$. For each state $s \in S$ and action $a \in A$, it yields the expected return if the agent starts in state $s$, takes action $a$, and then follows the policy for all future time steps. That is,
$$q_\pi(s,a) \doteq E_\pi[G_t | S_t = s, A_t = a]$$
We refer to $q_\pi(s,a)$ as the **value of taking action $a$ in state $s$ under a policy $\pi$** (or alternatively as the **value of the state-action pair $s,a$**).
- All optimal policies have the same action-value function $q_*$, called the **optimal action-value function**.

## Optimal Policies

---

- Once the agent determines the optimal action-value function $q_*$, it can quickly obtain an optimal policy $\pi_*$ by setting $\pi_*(s) = \mathrm{argmax}_{a \in A(s)} q_*(s,a)$.