where we just use the fact that we can express the value of the state-action pair $s_1, \text{right}$ as the sum of two quantities: (1) the immediate reward after moving right and landing on state $s_2$, and (2) the cumulative reward obtained if the agent begins in state $s_2$ and follows the policy.

Please now use the state-value function $v_\pi$ to calculate $q_\pi(s_1, \text{down})$, $q_\pi(s_2, \text{left})$, $q_\pi(s_2, \text{down})$, $q_\pi(s_3, \text{up})$, and $q_\pi(s_3, \text{right})$.

## For More Complex Environments

In this simple gridworld example, the environment is **deterministic**. In other words, after the agent selects an action, the next state and reward are 100% guaranteed and non-random. For deterministic environments, $p(s', r|s, a) \in \{0, 1\}$ for all $s', r, s, a$.

> In this case, when the agent is in state $s$ and takes action $a$, the next state $s'$ and reward $r$ can be predicted with certainty, and we must have $q_\pi(s, a) = r + \gamma v_\pi(s')$.

In general, the environment need not be deterministic, and instead may be **stochastic**. This is the default behavior of the `FrozenLake` environment from the mini project; in this case, once the agent selects an action, the next state and reward cannot be predicted with certainty and instead are random draws from a (conditional) probability distribution $p(s', r|s, a)$.

> In this case, when the agent is in state $s$ and takes action $a$, the probability of each possible next state $s'$ and reward $r$ is given by $p(s', r|s, a)$. In this case, we must have $q_\pi(s, a) = \sum_{s' \in \mathcal{S}^+, r \in \mathcal{R}} p(s', r|s, a)(r + \gamma v_\pi(s'))$, where we take the expected value of the sum $r + \gamma v_\pi(s')$.

Over the next couple concepts, you'll use this equation to write a function that yields an action-value function $q_\pi$ corresponding to a policy $\pi$ for the `FrozenLake` environment.