# MDPs: policy evaluation

# Evaluating a policy

**Definition: utility**

Following a policy yields a **random path**.

The **utility** of a policy is the (discounted) sum of the rewards on the path (this is a random variable).

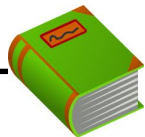| Path | Utility |
|------|---------|
| [in; stay, 4, end] | 4 |
| [in; stay, 4, in; stay, 4, in; stay, 4, end] | 12 |
| [in; stay, 4, in; stay, 4, end] | 8 |
| [in; stay, 4, in; stay, 4, in; stay, 4, in; stay, 4, end] | 16 |
| ... | ... |

**Definition: value (expected utility)**

The **value** of a policy at a state is the **expected** utility.
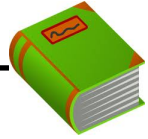
Value: 12

# Evaluating a policy: volcano crossing



Value: 3.73

Utility: -29.99

# Discounting

> **Definition: utility**
>
> Path: $s_0, a_1 r_1 s_1, a_2 r_2 s_2, \ldots$ (action, reward, new state).
> The **utility** with discount $\gamma$ is
> $$u_1 = r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 + \cdots$$

Discount $\gamma = 1$ (save for the future):

[stay, stay, stay, stay]: $4 + 4 + 4 + 4 = 16$
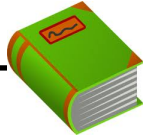
Discount $\gamma = 0$ (live in the moment):

[stay, stay, stay, stay]: $4 + 0 \cdot (4 + \cdots) = 4$

Discount $\gamma = 0.5$ (balanced life):

[stay, stay, stay, stay]: $4 + \frac{1}{2} \cdot 4 + \frac{1}{4} \cdot 4 + \frac{1}{8} \cdot 4 = 7.5$
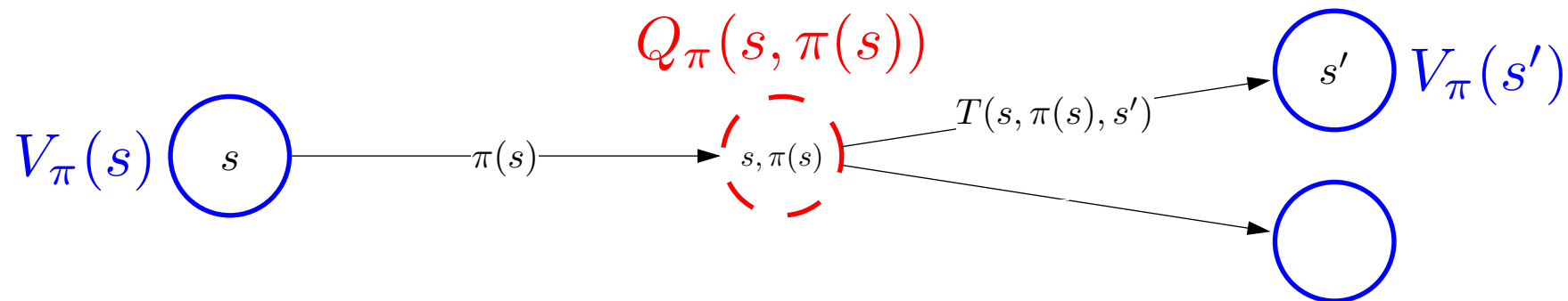
# Policy evaluation

**Definition: value of a policy**

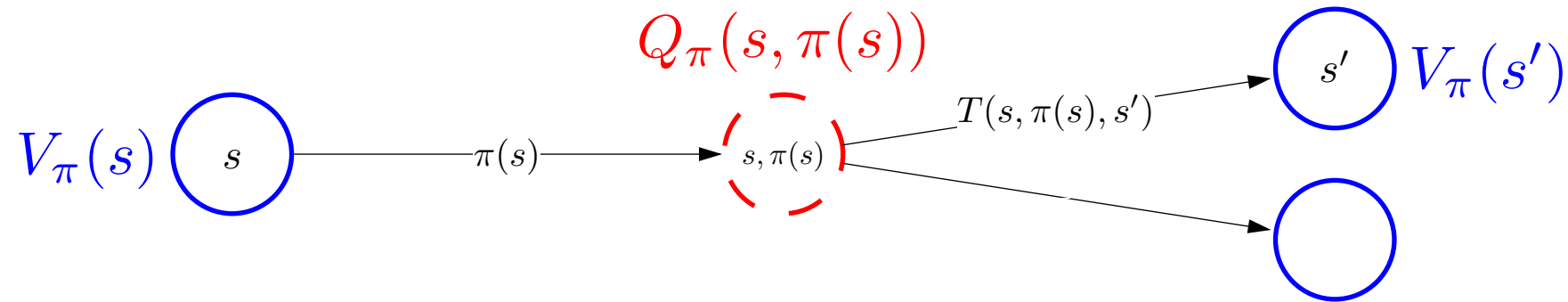Let $V_\pi(s)$ be the expected utility received by following policy $\pi$ from state $s$.

**Definition: Q-value of a policy**

Let $Q_\pi(s, a)$ be the expected utility of taking action $a$ from state $s$, and then following policy $\pi$.

# Policy evaluation
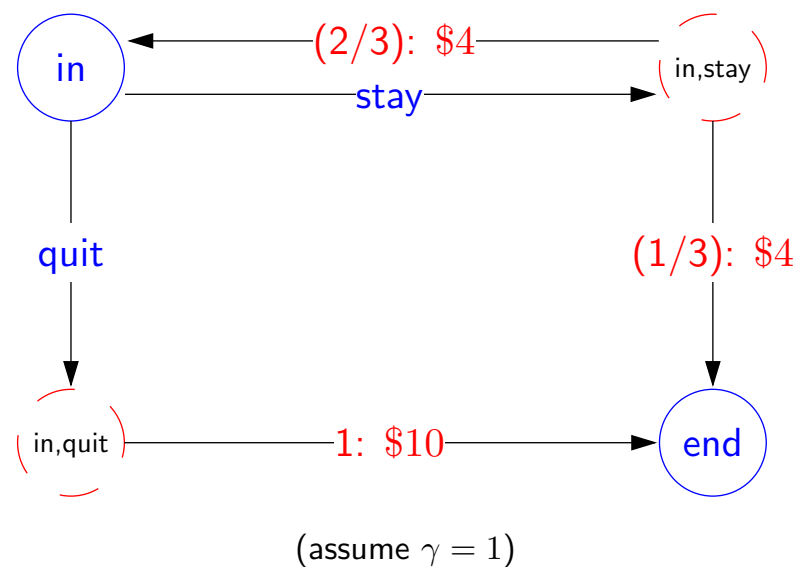
Plan: define recurrences relating value and Q-value



$$V_\pi(s) = \begin{cases} 0 & \text{if IsEnd}(s) \\ Q_\pi(s, \pi(s)) & \text{otherwise.} \end{cases}$$

$$Q_\pi(s, a) = \sum_{s'} T(s'|s, a)[\text{Reward}(s, a, s') + \gamma V_\pi(s')]$$

# Dice game



(assume $\gamma = 1$)

Let $\pi$ be the "stay" policy: $\pi(\text{in}) = \text{stay}$.

$$V_\pi(\text{end}) = 0$$

$$V_\pi(\text{in}) = \tfrac{1}{3}(4 + V_\pi(\text{end})) + \tfrac{2}{3}(4 + V_\pi(\text{in}))$$

In this case, can solve in closed form:

$$V_\pi(\text{in}) = 12$$

# Policy evaluation

**Key idea: iterative algorithm**

Start with arbitrary policy values and repeatedly apply recurrences to converge to true values.

**Algorithm: policy evaluation**

Initialize $V_\pi^{(0)}(s) \leftarrow 0$ for all states $s$.

For iteration $t = 1, \ldots, t_{\mathsf{PE}}$:

For each state $s$:

$$V_\pi^{(t)}(s) \leftarrow \underbrace{\sum_{s'} T(s'|s, \pi(s))[\mathsf{Reward}(s, \pi(s), s') + \gamma V_\pi^{(t-1)}(s')]}_{Q^{(t-1)}(s, \pi(s))}$$

# Policy evaluation implementation

How many iterations ($t_{\mathsf{PE}}$)? Repeat until values don't change much:

$$\max_{s \in \mathsf{States}} |V_\pi^{(t)}(s) - V_\pi^{(t-1)}(s)| \leq \epsilon$$

Don't store $V_\pi^{(t)}$ for each iteration $t$, need only last two:

$$V_\pi^{(t)} \text{ and } V_\pi^{(t-1)}$$

# Complexity

**Algorithm: policy evaluation**

Initialize $V_\pi^{(0)}(s) \leftarrow 0$ for all states $s$.

For iteration $t = 1, \ldots, t_{\mathsf{PE}}$:

    For each state $s$:

$$V_\pi^{(t)}(s) \leftarrow \underbrace{\sum_{s'} T(s'|s, \pi(s))[\mathsf{Reward}(s, \pi(s), s') + \gamma V_\pi^{(t-1)}(s')]}_{Q^{(t-1)}(s, \pi(s))}$$

**MDP complexity**

$S$ states

$A$ actions per state

$S'$ successors (number of $s'$ with $T(s'|s, a) > 0$)

Time: $O(t_{\mathsf{PE}} S S')$

# Policy evaluation on dice game

Let $\pi$ be the "stay" policy: $\pi(\text{in}) = \text{stay}$.

$$V_\pi^{(t)}(\text{end}) = 0$$

$$V_\pi^{(t)}(\text{in}) = \tfrac{1}{3}(4 + V_\pi^{(t-1)}(\text{end})) + \tfrac{2}{3}(4 + V_\pi^{(t-1)}(\text{in}))$$

| $s$ | end | in | |
|---|---|---|---|
| $V_\pi^{(t)}$ | 0.00 | 12.00 | $(t = 100 \text{ iterations})$ |

Converges to $V_\pi(\text{in}) = 12$.

# Summary so far

- **MDP**: graph with states, chance nodes, transition probabilities, rewards

- **Policy**: mapping from state to action (solution to MDP)

- **Value of policy**: expected utility over random paths

- **Policy evaluation**: iterative algorithm to compute value of policy