# Section B

The goal of RL is to find an __optimal policy__:

$$\pi^* = \underset{\pi}{\text{argmax}} \ V^{\pi}(s) \quad \forall s \in S$$

__Value function__:

$$V^{\pi}(s_0) = \mathbb{E}\left[\text{total discounted reward for starting in state } s_0 \text{ and executing } \pi\right]$$

Given $\pi$, $p(s_{t+1}|s_t, a_t)$ there exists a distribution over __state trajectories__

$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \xrightarrow{a_3} \cdots \quad (\text{because } p(s'|s,a) \text{ is stochastic})$$

$$= \mathbb{E}_{p(s'|s,a)}\left[\underbrace{R(s_0,a_0)}_{r_0} + \underbrace{\gamma R(s_1,a_1)}_{\gamma r_1} + \underbrace{\gamma^2 R(s_2,a_2)}_{\gamma^2 r_2} + \cdots \mid s_0\right]$$

$$= R(s_0,a_0) + \gamma \ \mathbb{E}_{p(s'|s,a)}\left[\underbrace{R(s_1,a_1)}_{f(s_1)} + \underbrace{\gamma R(s_2,a_2) + \gamma^2 R(s_3,a_3) + \cdots}_{g(s_2,s_3,\cdots)} \mid s_0\right]$$
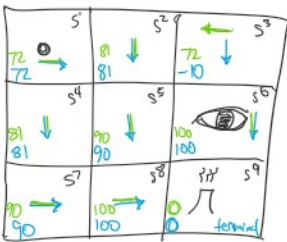
$$= R(s_0,a_0) + \gamma \sum_{s_1 \in S} p(s_1|s_0,a_0)\left(R(s_1,a_1) + \gamma \ \mathbb{E}_{p(s'|s,a)}\left[R(s_2,a_2) + \gamma R(s_3,a_3) + \cdots \mid s_0\right]\right)$$

$$\underbrace{\phantom{V^{\pi}(s_1)}}_{V^{\pi}(s_1)}$$

$$= R(s_0,a_0) + \gamma \sum_{s_1 \in S} p(s_1|s_0,a_0) V^{\pi}(s_1)$$

### Bellman Equations:

$$V^{\pi}(s) = R(s,\pi(s)) + \gamma \sum_{s' \in S} p(s'|s,\pi(s)) V^{\pi}(s')$$

| Ex: Value functions and policies |



Assume $\gamma = 0.9$

$A = \{\leftarrow, \downarrow, \uparrow, \rightarrow\}$

$R(s,a) = +100$ if entering 🏆

$R(s,a) = -100$ if entering 👁

$R(s,a) = 0$

Transitions are deterministic, $\delta(s,a)$

$$\begin{cases} V^{\pi}(s) = R(s,\pi(s)) + \gamma V^{\pi}(\delta(s,\pi(s))) \\ \leftarrow \text{some policy } \pi' \\ V^{\pi'}(s^3) = -100 + (0.9)100 = -10 \end{cases}$$

$$\begin{cases} V^*(s) = \max_{a \in A} R(s,a) + \gamma V^*(\delta(s,a)) \\ \leftarrow \text{optimal policy } \pi^* \\ V^*(s^3) = \max\left( \begin{matrix} R(s^3,\downarrow) + 0.9 \ V^*(s^6) \\ R(s^3,\leftarrow) + 0.9 \ V^*(s^2) \end{matrix} \right), = 72 \end{cases}$$

- __Optimal value function__  $V^* = V^{\pi^*}$

- Given $V^*, R(s,a), p(s'|s,a), \gamma$ we can compute $\pi^*$!

$$\pi^*(s) = \underset{a \in A}{\text{argmax}} \ \underbrace{R(s,a)}_{\substack{\text{immediate} \\ \text{reward}}} + \underbrace{\gamma \sum_{s' \in S} p(s'|s,a) V^*(s')}_{\substack{\text{future discounted} \\ \text{reward}}}$$

☹ Problem: our definitions of $\pi^*$ and $V^*$ are cyclic ☹

- Can compute $V^*$ without $\pi^*$

$$V^*(s) = \max_{a \in A} R(s,a) + \gamma \sum_{s' \in S} p(s'|s,a) V^*(s')$$

$\leftarrow$ ..... definition of optimal value function

$$V^*(s) = \max_{a \in A} R(s,a) + \sum_{s' \in S} p(s' \mid s, a) + (-)$$

★ recursive definition of optimal value function
system of $|S|$ equations and $|S|$ variables
each variable is $V^*(s)$ for some $s \in S$

Key Idea of value iteration:

- Apply dynamic programming, i.e. fixed point iteration, to the recursive def. of $V^*$