# Reinforcement learning

Leonardo Toffalini

2026-02-21

# Outline

# 1. Informal introduction to RL

What makes RL different?

- Not supervised (learning from labeled data)
- Not unsupervised (learning patterns in unlabeled data)
- Only *reward* signal
- Feedback may be delayed
- Heavily time dependent
- The agent has affect on the data

# 1.2 Examples of RL

- Robotics
- Video games
- Self-driving
- Finance
- Natural language processing (recently)
- Recommendation systems
- Many more...

# 2. Setting the scene

**Definition 2.1.1** A Markov decision process (MDP) is defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}_a, \mathcal{R}_a)$, where
- $\mathcal{S}$ is the set of all states
- $\mathcal{A}$ is the set of all possible actions
- $\mathcal{P}$ is the transition probability function

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' \mid S_t = s, \quad A_t = a]$$

- $\mathcal{R}$ is the reward function

$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, \quad A_t = a]$$

**Definition 2.1.2** A state $S_t$ is Markov if and only if

$$\mathbb{P}[S_{t+1} \mid S_t] = \mathbb{P}[S_{t+1} \mid S_1, ..., S_t].$$

## 🗨 Remark

In a Markov decision process, the successor state is *solely* influenced by the current state.

That is, an MDP has no memory of previous states.
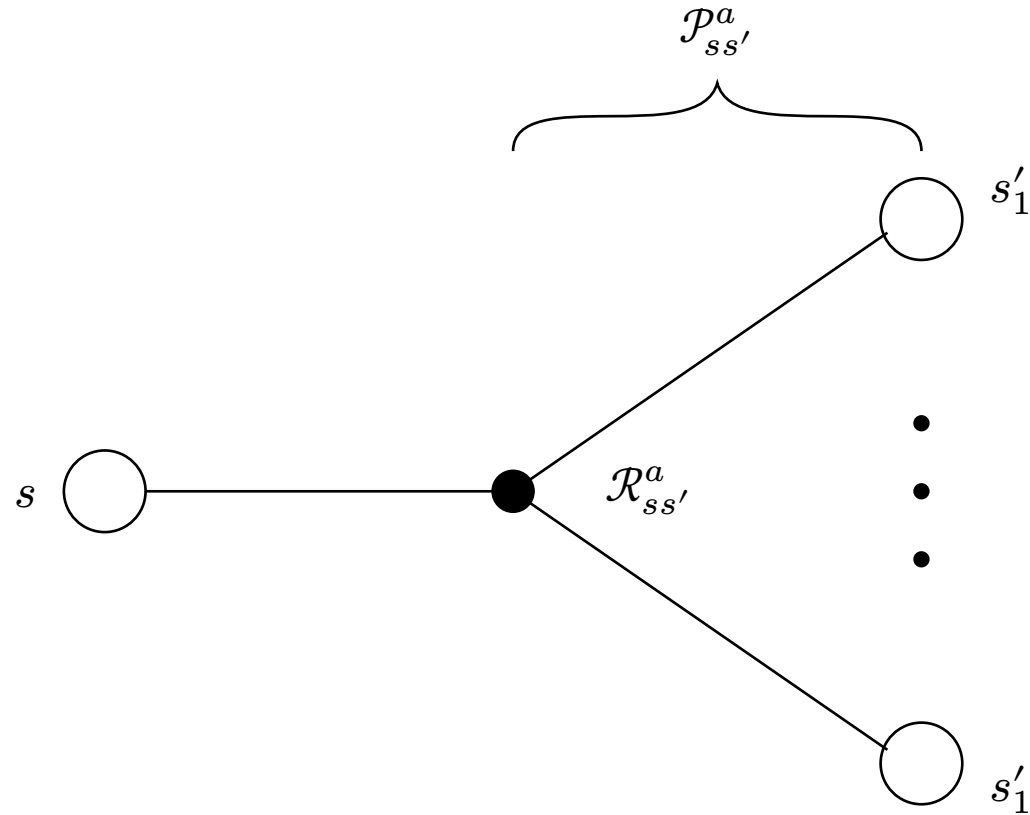
Is this a limitation?

Figure 1: hollow circles = states, full circles = actions

**Definition 2.3.1** The reward hypothesis claims that all goals can be described by the maximization of expected cumulative reward.

## 💬 **Remark**

Do you agree?

Can you find a counterexample?

**Definition 2.4.1** A policy $\pi$ is a distribution over the action space conditioned on the state space, that is

$$\pi(a|s) = \mathbb{P}[A_t = a \mid S_t = s].$$

## 💬 Remark

The next action depends *only* on the current state and nothing else.

The policy may be deterministic if $\pi(a_i \mid s) = 1$ for some $i$ and $\pi(a_j \mid s) = 0$ for any other $j \neq i$.

**Definition 2.5.1** The return $G_t$ is the total discounted reward from time step $t$, that is

$$G_t = R_{t+1} + \gamma R_{t+2} + ... = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}.$$

💬 **Remark**

We can control the importance of *immediate* versus *future* rewards with $\gamma \in [0, 1]$.

**Definition 2.5.2** The state-value function of a state $s$ is the expected discounted return starting from state $s$, that is

$$v_\pi(s) := \mathbb{E}_\pi[G_t \mid S_t = s].$$

**Definition 2.5.3** The action-value function of a state-action pair $(s, a)$ is the expected discounted return starting from state $s$ and taking action $a$, that is

$$q_\pi(s, a) := \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

# 3. Bellman equations

> **Proposition 3.1.1**
>
> $$v(s) = \mathbb{E}_\pi\big[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s\big]$$

*Proof.*

$$
\begin{aligned}
v_\pi(s) &:= \mathbb{E}_\pi[G_t \mid S_t = s] \\
&= \mathbb{E}_\pi\big[R_{t+1} + \gamma G_{t+1} \mid S_t = s\big] \\
&= \mathbb{E}_\pi\big[R_{t+1} + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_{t+1}] \mid S_t = s\big] \\
&= \mathbb{E}_\pi\big[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s\big]
\end{aligned}
$$

$\square$

**Proposition 3.1.2**

$$q_\pi(s, a) = \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

*Proof.*

$$
\begin{aligned}
q_\pi(s, a) &:= \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_{t+1}, A_{t+1}] \mid S_t = s, A_t = a] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]
\end{aligned}
$$

$\square$

**Proposition 3.1.3**

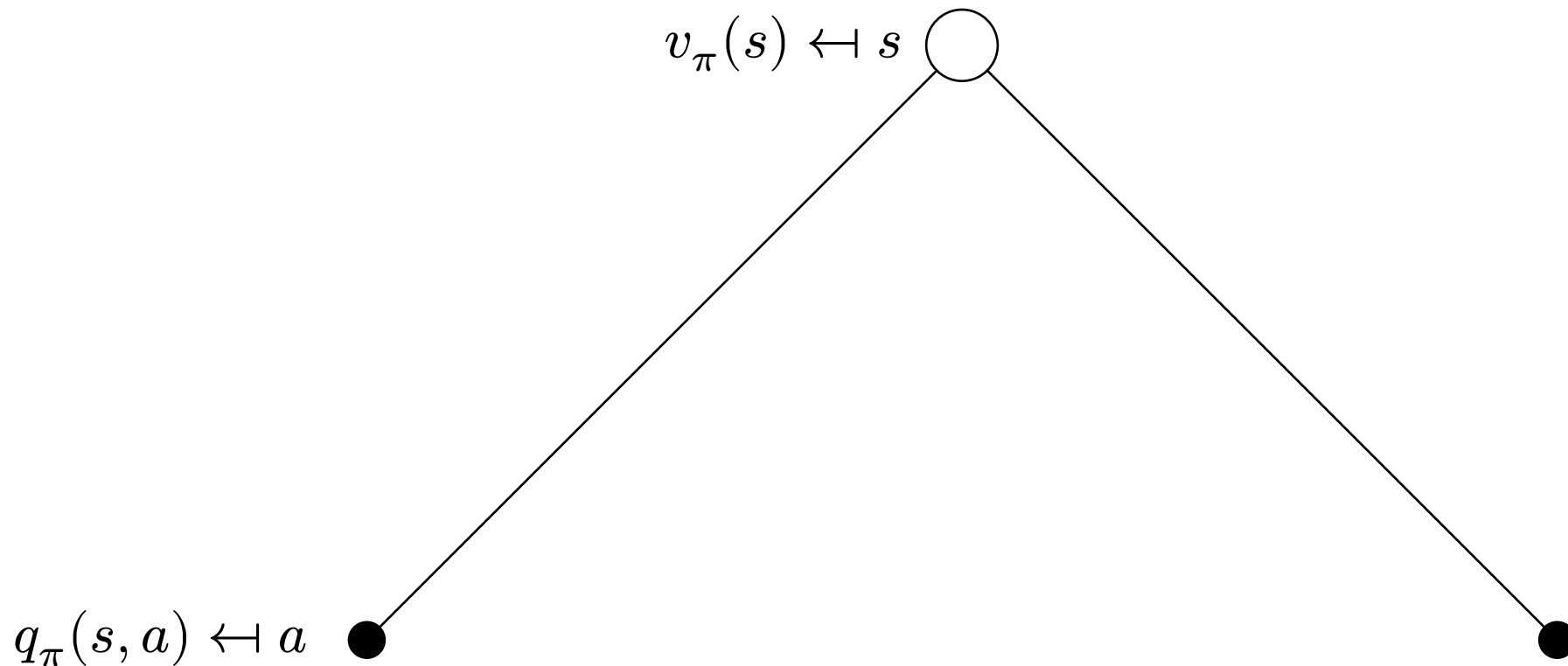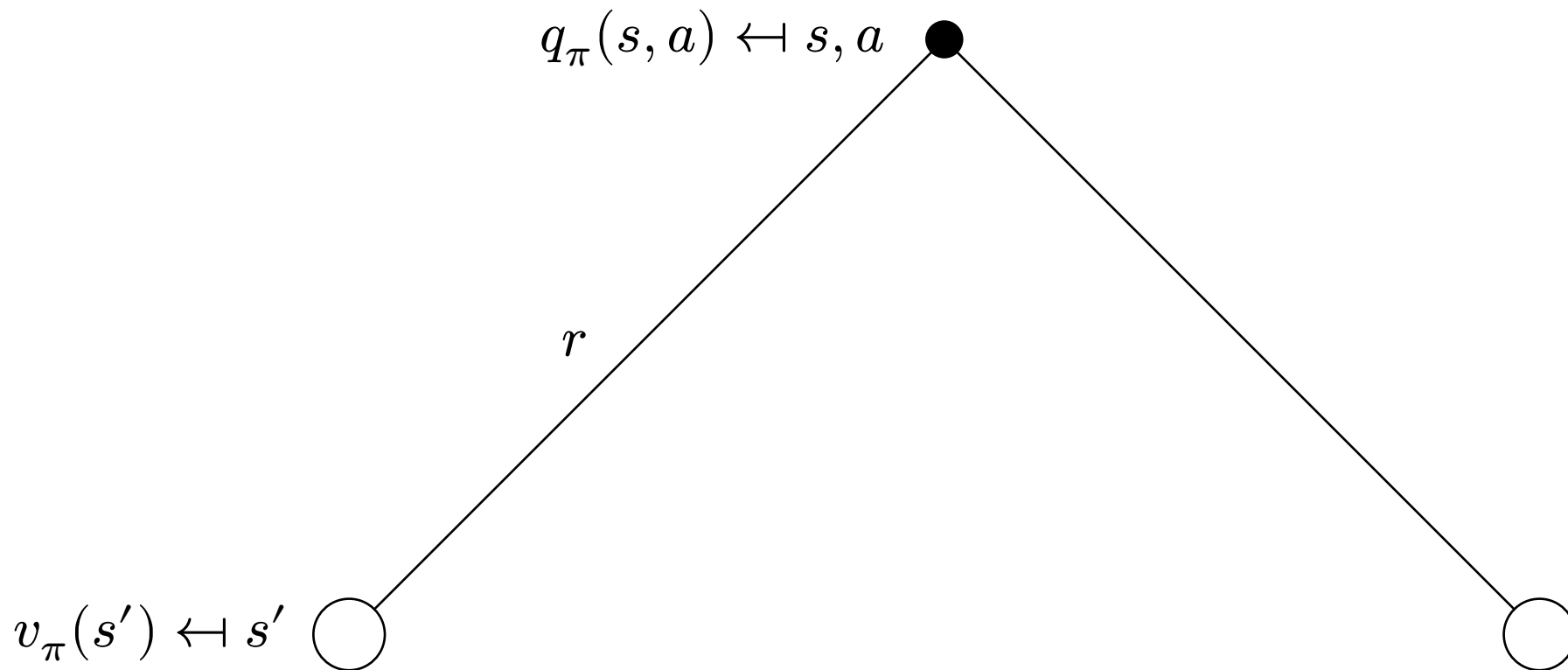$$v_\pi(s) = \sum_{a\in\mathcal{A}} \pi(a|s)q_\pi(s,a)$$

*Proof.*

$$v_\pi(s) := \mathbb{E}_\pi[G_t \mid S_t = s]$$
$$= \mathbb{E}_\pi[\mathbb{E}_\pi[G_t \mid S_t, A_t] \mid S_t = s]$$
$$= \sum_{a\in\mathcal{A}} \pi(a|s)q_\pi(s,a)$$

$\square$

$$v_\pi(s) \leftarrowtail s \quad \bigcirc$$

$$q_\pi(s, a) \leftarrowtail a \quad \bullet$$

**Proposition 3.3.1**

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$

*Proof.*

$$q_\pi(s, a) := \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[R_t + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[R_t + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_{t+1} = s] \mid S_t = s, A_t = a]$$

$$= \mathbb{E}_\pi[R_t + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \mathcal{R}_s^a + \gamma \mathbb{E}_\pi[v_\pi(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$

$\square$

$$q_\pi(s, a) \longleftarrow s, a$$

$$r$$

$$v_\pi(s') \longleftarrow s'$$

**Proposition 3.5.1**
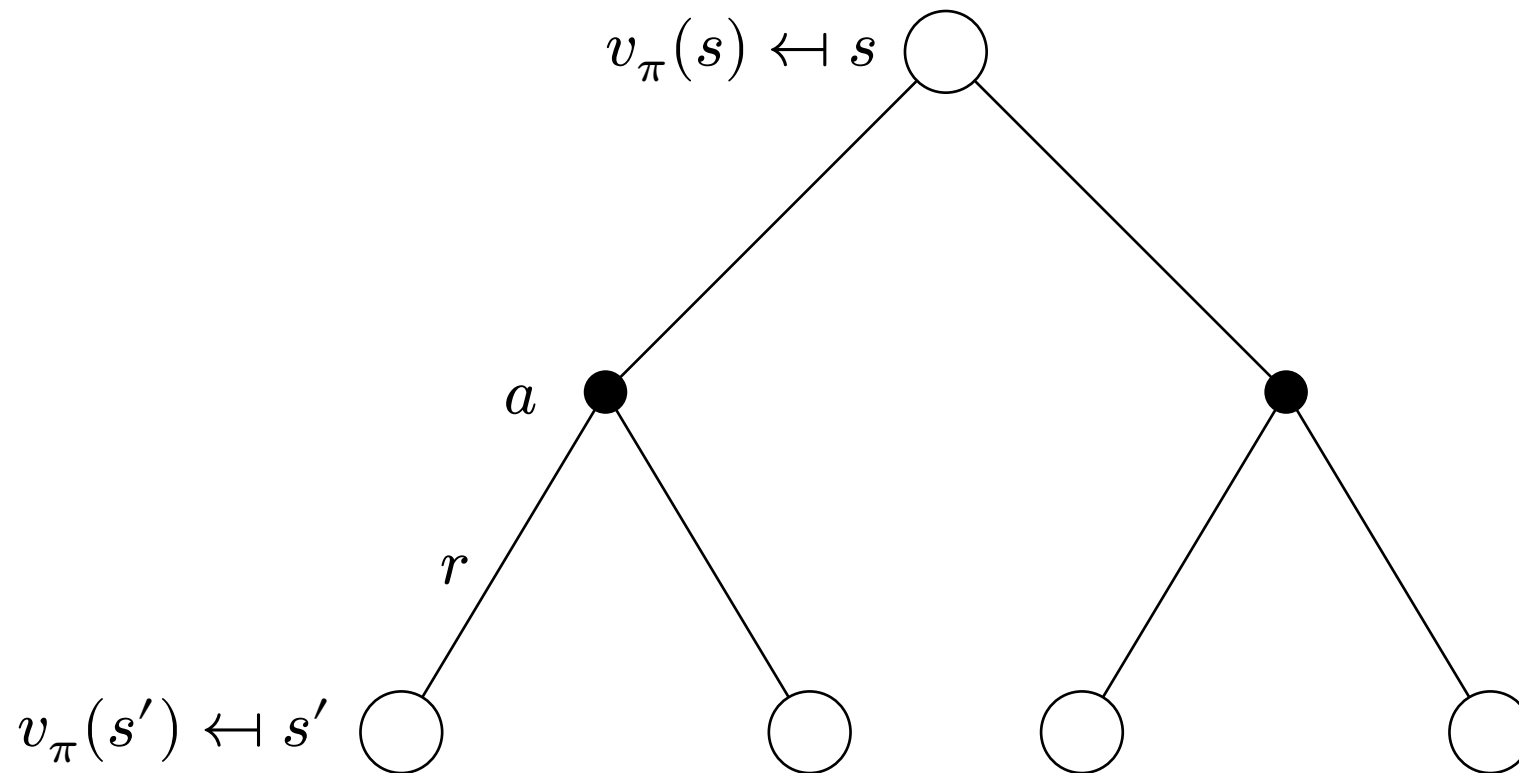
$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a v_\pi(s') \right)$$

**Proposition 3.5.2**

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') q_\pi(s', a')$$
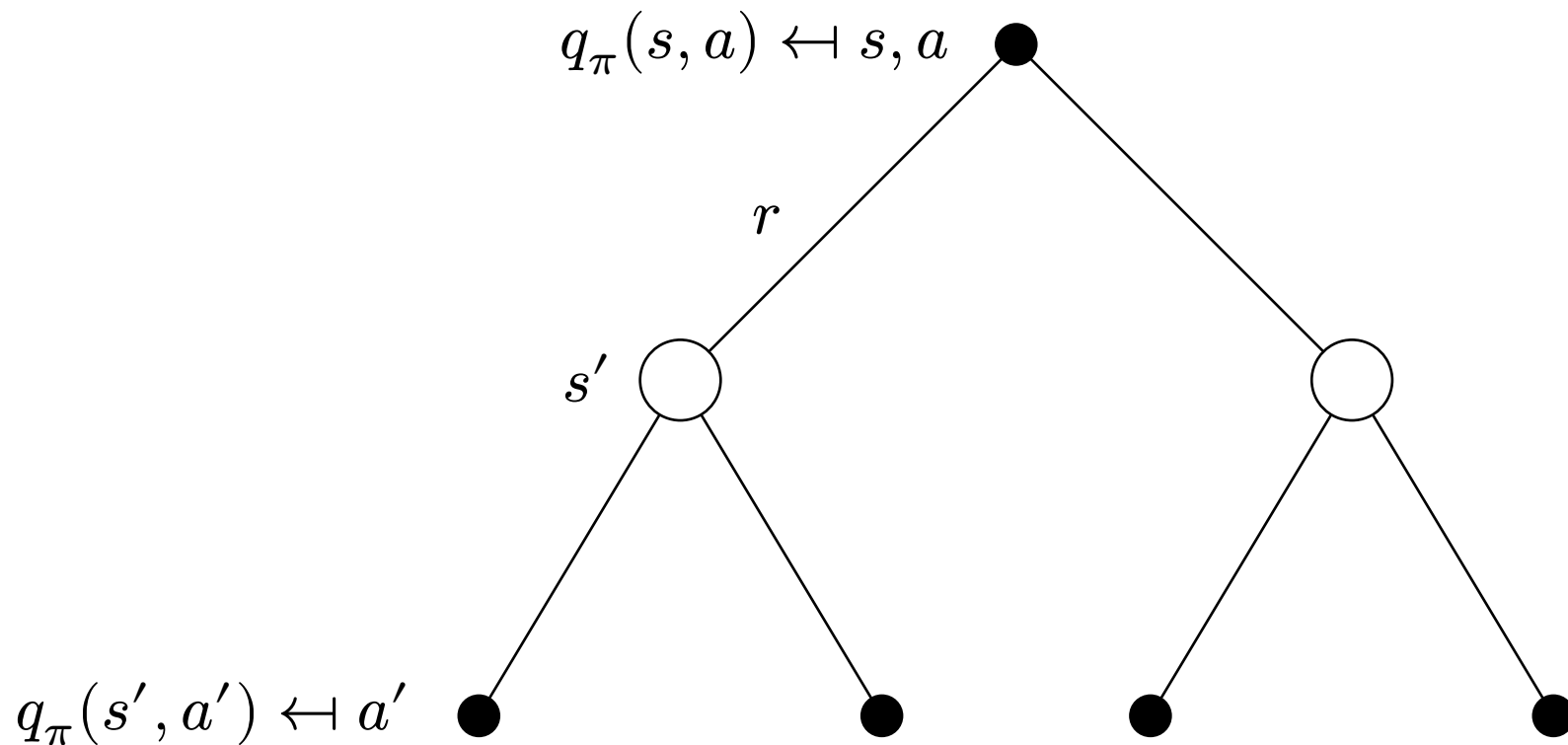
$q_\pi(s, a) \leftarrowtail s, a$

$r$

$s'$

$q_\pi(s', a') \leftarrowtail a'$

**Definition 3.7.1** The optimal state-value function is defined as

$$v_*(s) = \max_\pi v_\pi(s)$$

**Definition 3.7.2** The optimal action-value function is defined as

$$q_*(s,a) = \max_\pi q_\pi(s,a).$$

**Definition 3.7.3**  A policy $\pi^*$ is optimal if for all other policies $\pi$

$$v_{\pi^*}(s) \geq v_\pi(s) \quad \forall s \in \mathcal{S}.$$

> ### Remark

That is, $\pi^*$ is optimal if the value function induced by following $\pi^*$ is optimal.
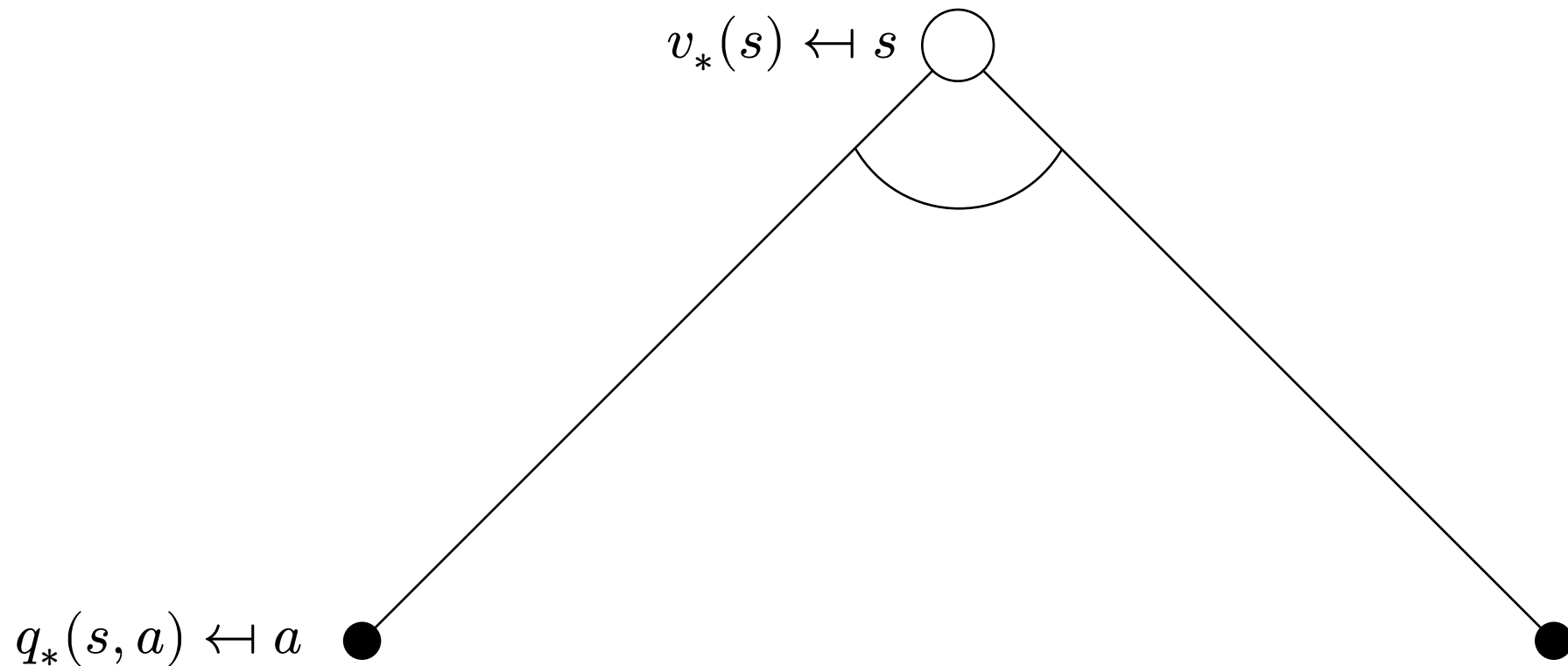
**Proposition 3.7.4**

$$v_*(s) = \max_a q_*(s, a)$$

**Proposition 3.7.5**

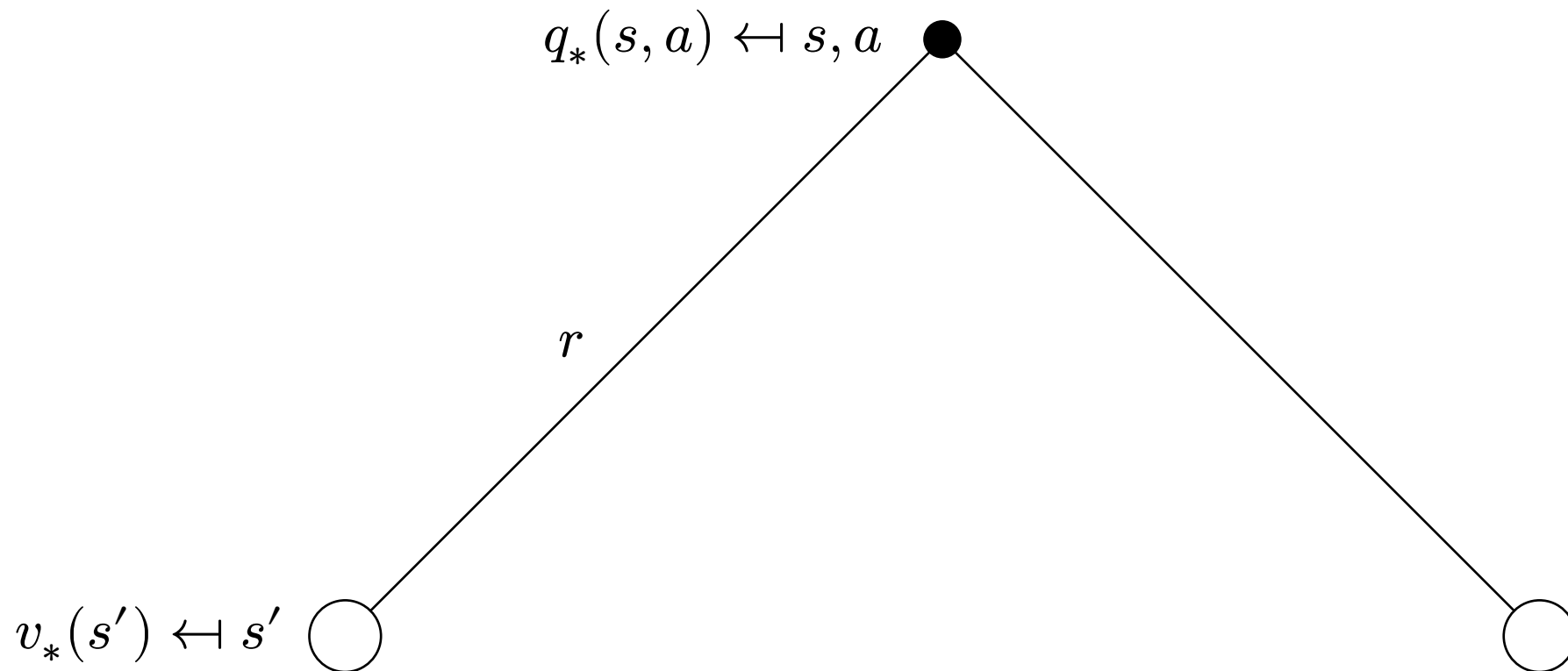$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

$$v_*(s) \;\longleftmapsto\; s \;\bigcirc$$

$$q_*(s,a) \;\longleftmapsto\; a \;\bullet$$

$q_*(s,a) \leftmapsto s,a$
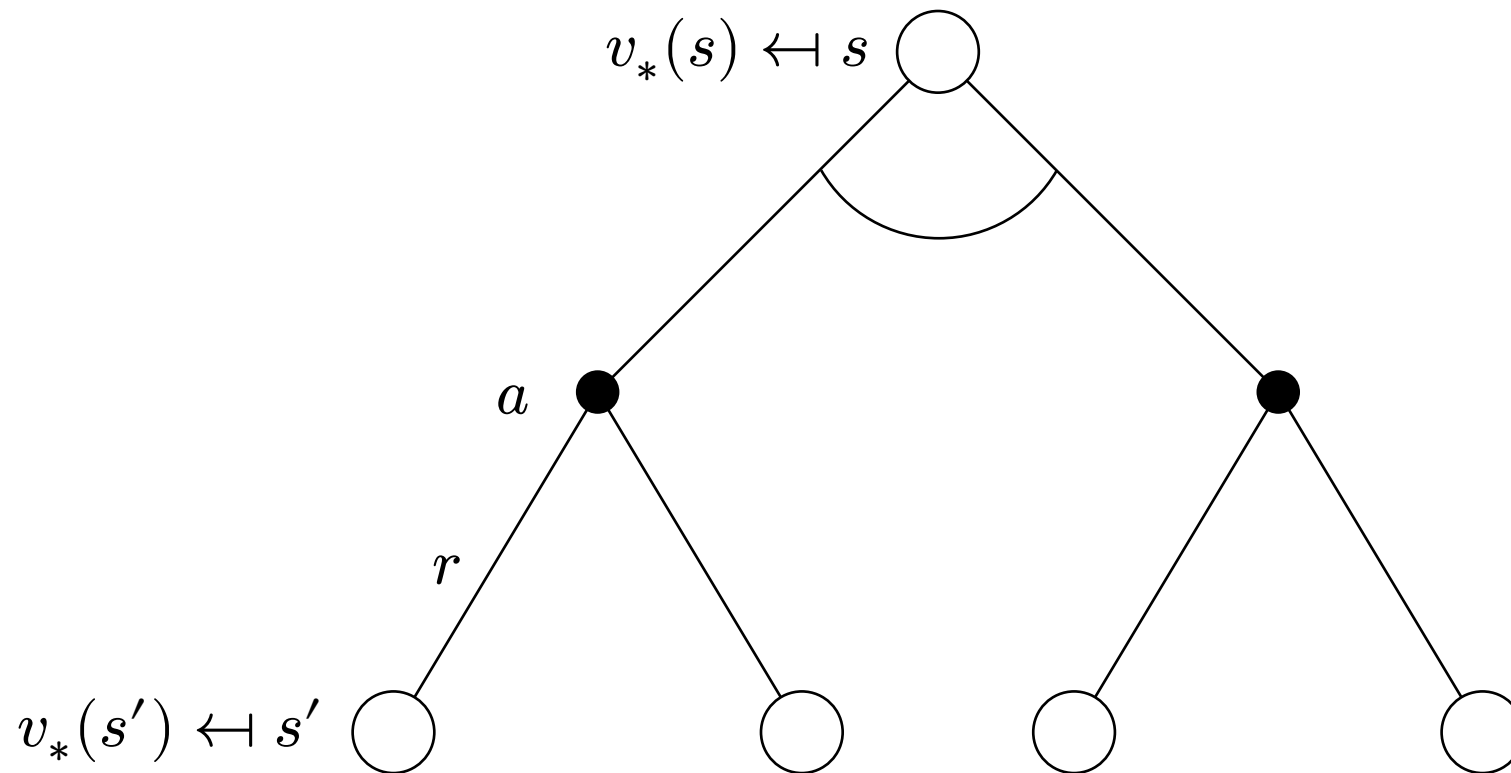
$r$

$v_*(s') \leftmapsto s'$

**Proposition 3.9.1**

$$v_*(s) = \max_a \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \right)$$

**Proposition 3.9.2**

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a'} q_*(s', a')$$
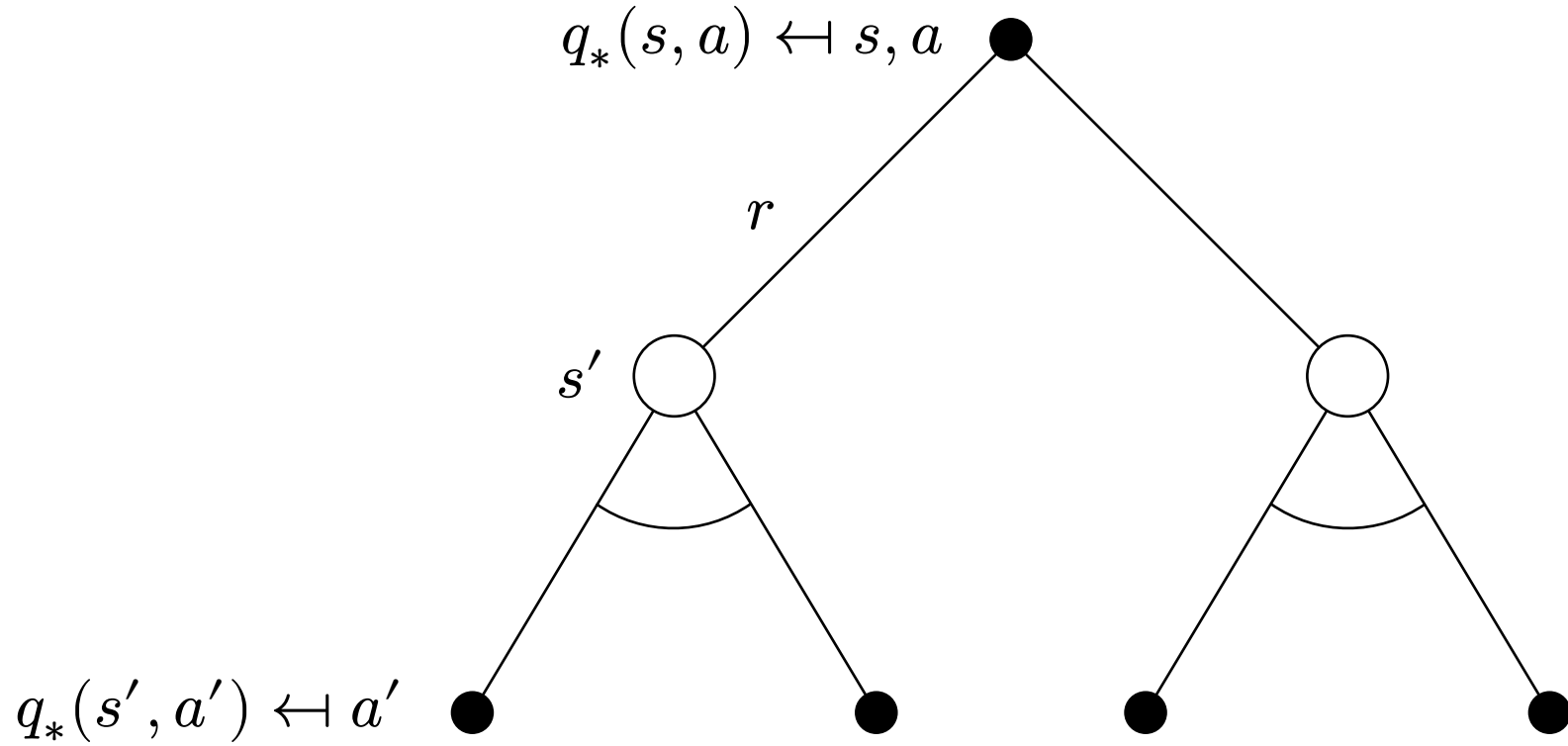
Now lets see what we can do with this theory.

Lets see some algorithms in the practice session...