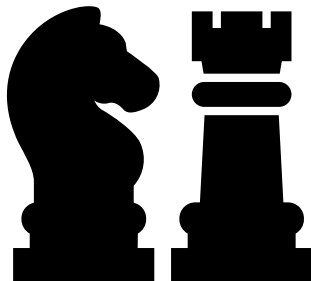# Reinforcement Learning 01

DS Development Presentations
Peter Nicholas S. Onglao | MNL

**GOAL:**

Make an RL algorithm for chess (>3 presentations away) that can beat me

# 📌 Contents

Re-Introduction to Reinforcement Learning (RL)

Markov Decision Processes (MDP)

- <u>UCL Course on Reinforcement Leaning</u>, David Silver

  - ➤ Video lectures with accompanying slides

  - ➤ David Silver: AlphaGo, AlphaZero, AlphaStar (DeepMind)

  - ➤ Content is patterned after his lectures

- <u>Reinforcement Learning – An Introduction</u>, Sutton and Barto

# 1

## What is RL?

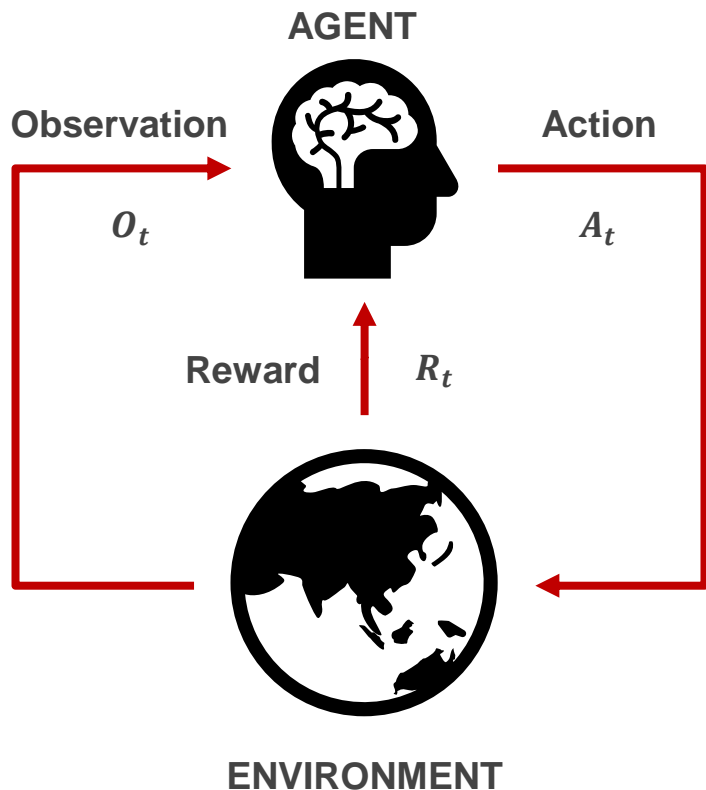**Introduction + basic math requirements for an RL problem.**

What makes RL different from supervised/unsupervised learning?
- No supervisor, only a **reward**
- Feedback is delayed, not instantaneous
- **Time really matters**, not i.i.d. data
- The agent affects the data it receives

Basic components:
- Agent
- Environment
- Observation
- Action
- Reward

**AGENT**

**Observation**

$O_t$

**Action**

$A_t$

**Reward** $R_t$

**ENVIRONMENT**

At each step $t$ the agent:
- ▷ Executes action $A_t$
- ▷ Receives observation $O_t$
- ▷ Receives reward $R_t$

At each step $t$ the environment:
- ▷ Receives action $A_t$
- ▷ Emits observation $O_{t+1}$
- ▷ Emits reward $R_{t+1}$

# Reward

- The reward $R_t$ is a **SCALAR** feedback signal
- Indicates how well the agent is doing at step $t$
- The agent's job is to maximize **cumulative reward**
  - Cumulative reward since actions may have long-term consequences, and reward may be delayed
  - Immediate vs. long-term gain
- **Reward hypothesis**
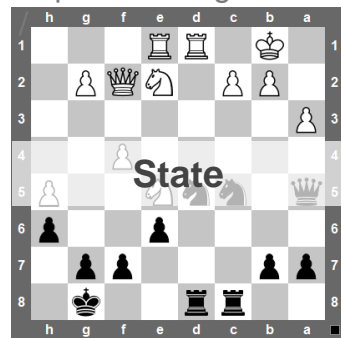  - All goals can be described by the maximization of expected cumulative reward

How will the agent act? History and State

**History**: the sequence of observations, actions, rewards

$$H_t = O_1, R_1, A_1, \ldots, A_{t-1}, O_t, R_t$$

All observable variables up to time $t$

The future depends on the history

- Agent: selects action
- Environment: selects observation/reward

**State**: a **concise** way to represent history

- It is a function of history
- $S_t = f(H_t)$

**History**

https://mat3e.github.io

**State**

chessfox.com

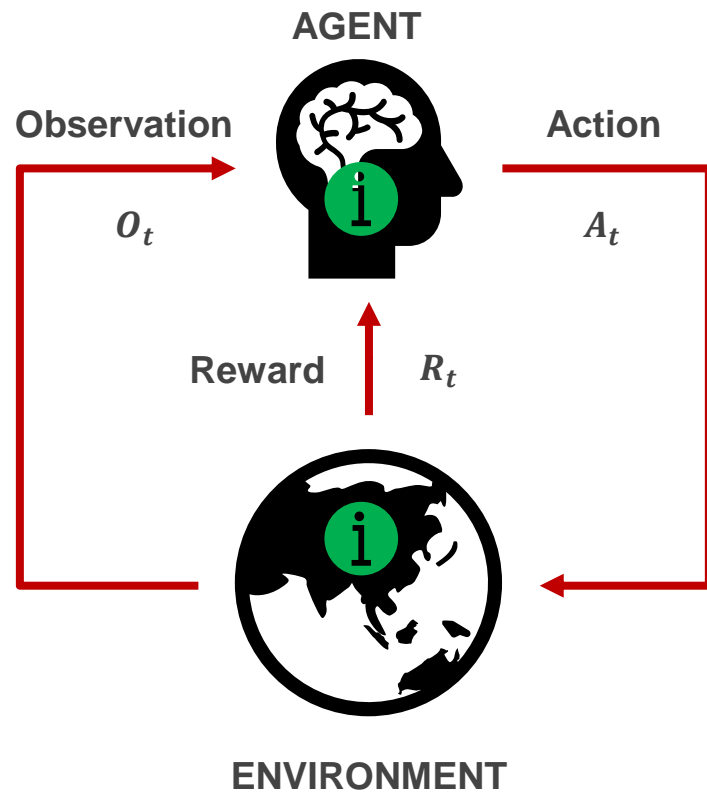- **Environment State** $(S_t^e)$
  - ➤ The **internal state** of the environment (some set of numbers)
  - ➤ The data the environment uses to pick the next observation/reward
  - ➤ Not usually visible to agent
  - ➤ If it is visible, it may not be useful to the agent
- **Agent State** $(S_t^a)$
  - ➤ The set of numbers inside the algorithm
  - ➤ The information used by the agent to pick the next action
  - ➤ **Our role** – choose which information to keep/take

**AGENT**

Observation

Action

$O_t$

$A_t$

Reward

$R_t$

**ENVIRONMENT**

**Markov State**
- ▷ Contains all useful information from the history
- ▷ A state $S_t$ is Markov iff:

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

"The future is independent of the past given the present"
- ▷ $H_{1:t} \to S_t \to H_{t+1:\infty}$

Once the current state is known, the **history may be thrown away**

The state is a sufficient statistic of the future

Environment state $S_t^e$ is Markov

History $H_t$ is Markov

# Components of an RL Agent

An RL agent may include **one or more** of the following:
- **Policy**: agent's behavior function
- **Value function**: how good is each state and/or action
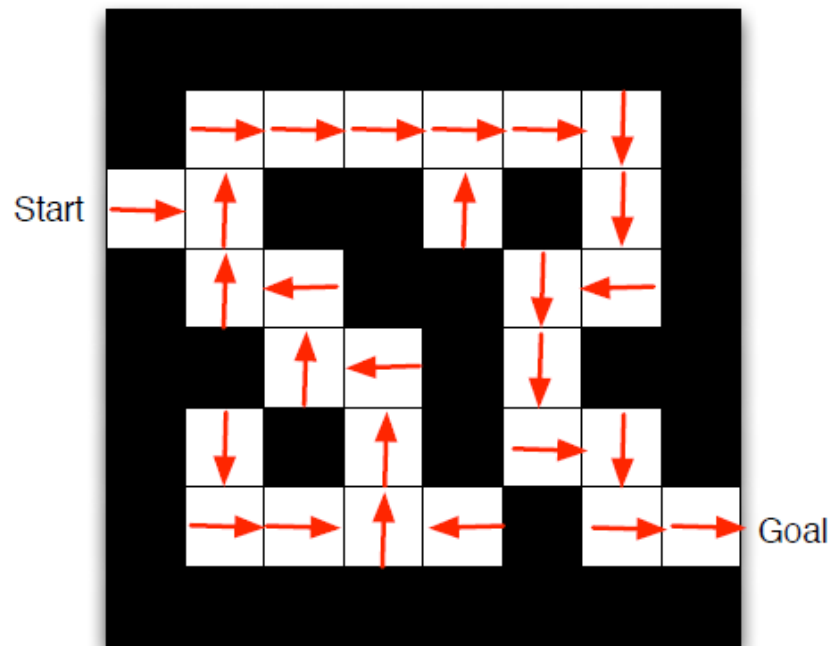- **Model**: agent's representation of the environment

- **Policy**: the agent's behavior; a map from state to action
  - ➤ Deterministic: $a = \pi(s)$
  - ➤ Stochastic: $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$
- **Value function**: prediction of future reward; evaluate the current state
  - ➤ $v_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+2} + \cdots | S_t = s]$
- **Model**: predicts what the environment will do next
  - ➤ $\mathcal{P}$ predicts the **next state**
    - ➤ $\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
  - ➤ $\mathcal{R}$ predicts the **next reward**
    - ➤ $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$

Credits: David Silver's slides
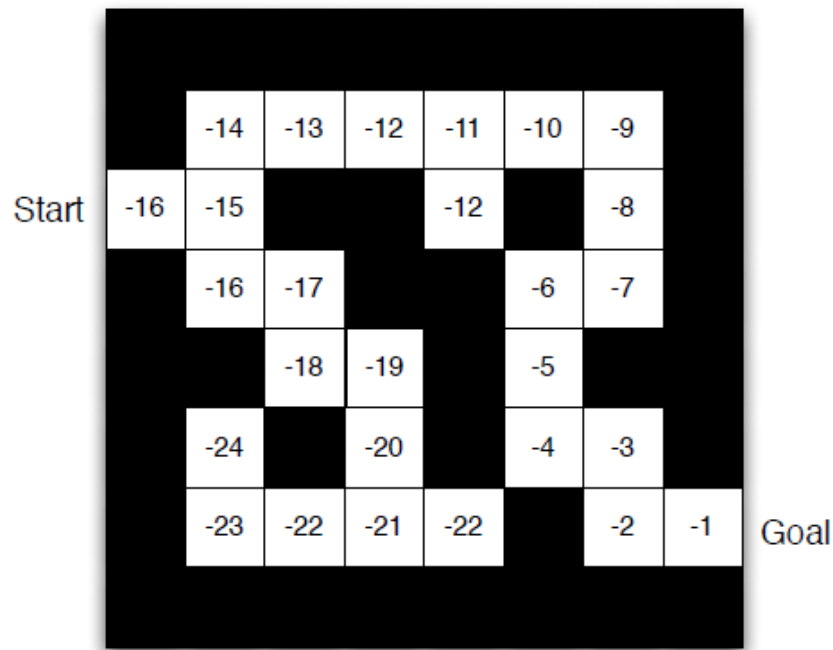
Each state $s$ = a position in the maze
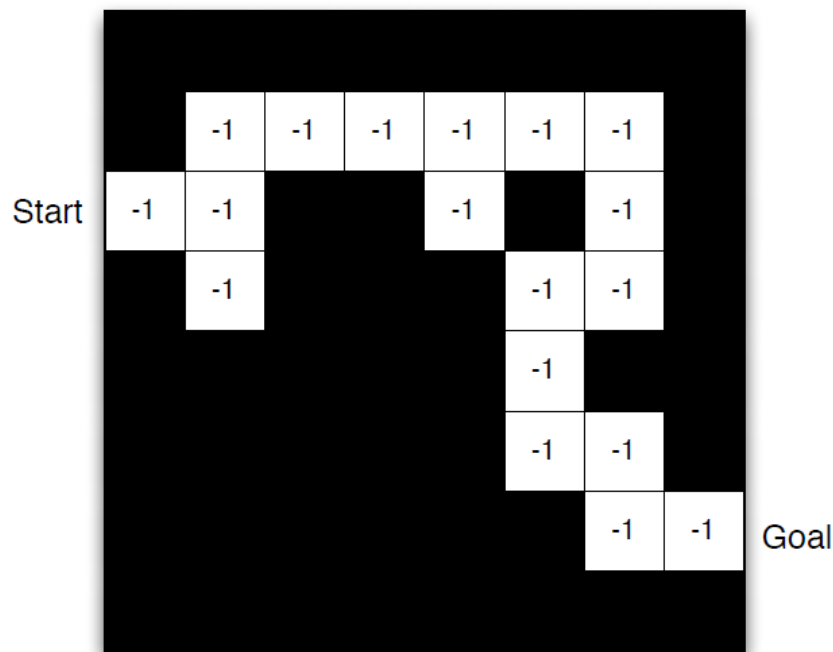
Arrows represent policy $\pi(s)$ for each state $s$

Credits: David Silver's slides

- Each state $s$ = a position in the maze
- Reward: -1 per step
- Numbers represent the value function $v_\pi(s)$ of each state $s$

Start

Goal

Credits: David Silver's slides

- Agent may have an **internal model** of the environment
- Model may be imperfect
- Grid layout represents transition model $\mathcal{P}_{ss'}^a$
- Numbers represent immediate reward $R_s^a$ from each state $s$
  - ▷ All the same since the *immediate reward* is the same each step

17

**(1) Depending on Policy & Value Function**
- ➤ Value Based:     No Policy (Implicit) + Value Function
- ➤ Policy Based:     Policy + No Value Function
- ➤ Actor Critic:     Policy + Value Function

**(2) Depending on Model**
- ➤ Model Free:     Policy and/or Value Function + No Model
- ➤ Model Based:     Policy and/or Value Function + Model

- Learning vs. Planning
  - **Learning**
    - **Environment** is initially **unknown**
    - Agents interacts w/ environment and improves its policy
  - **Planning**
    - A model of the **environment** is **known**
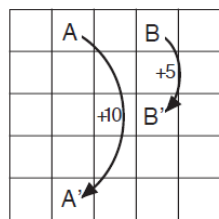    - Agent does not perform interactions but performs computations, and improves its policy this way

- Exploration vs Exploitation
  - ▷ **Exploration**: Discover a good policy
  - ▷ **Exploitation**: W/o losing too much reward along the way
- Prediction vs Control
  - ▷ **Prediction**: evaluate the future given a policy
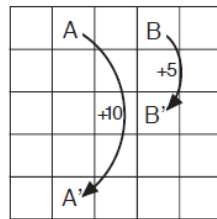  - ▷ **Control**: optimize the future and find the best policy



**Prediction** (policy = random)

**Control**

Credits: David Silver's slides

**Components of an RL problem**
- ➤ Agent, Environment
- ➤ Observation, Reward, Action
- ➤ History and State

**Components of an RL agent**
- ➤ Policy, Value Function, Model
- ➤ Agents can be categorized by presence/absence of those three

**Differentiating some RL terms**
- ➤ Learning vs Planning
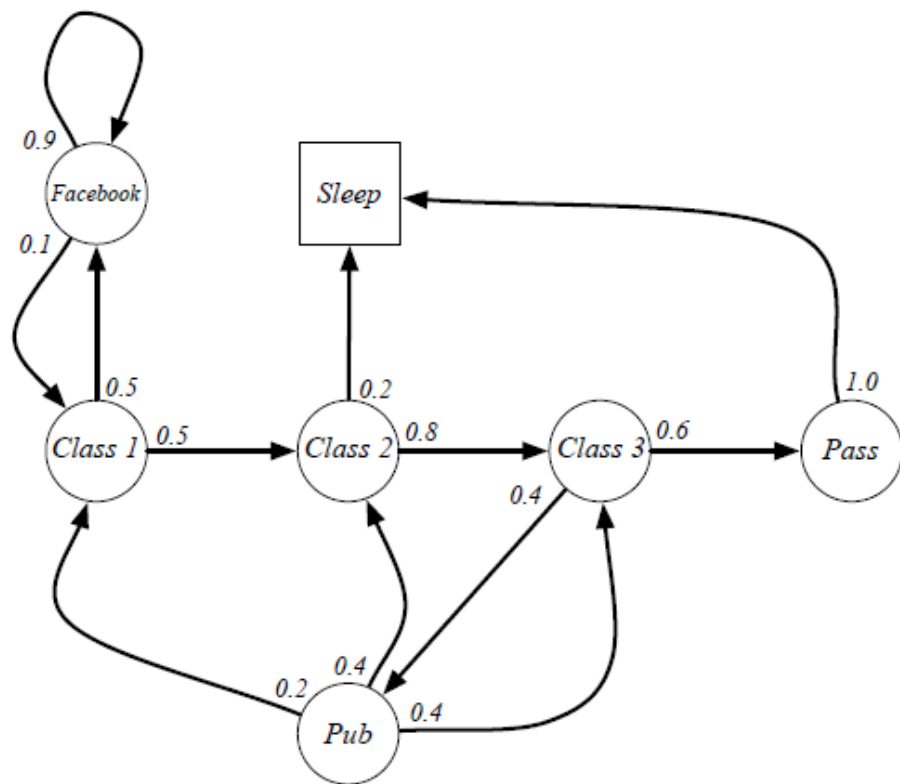- ➤ Exploration vs Exploitation
- ➤ Prediction vs Control

# 2

## Markov Decision Process

MDPs, MRPs, and the Bellman Equation

**Markov Process**/Markov Chain

➤ A tuple $\langle \mathcal{S}, \mathcal{P} \rangle$

➤ $\mathcal{S}$ is a (finite) set of states

➤ $\mathcal{P}$ is a state transition probability matrix

$$\mathcal{P} = \begin{array}{c|ccccccc} & C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ \hline C1 & & 0.5 & & & & 0.5 & \\ C2 & & & 0.8 & & & & 0.2 \\ C3 & & & & 0.6 & 0.4 & & \\ Pass & & & & & & & 1.0 \\ Pub & 0.2 & 0.4 & 0.4 & & & & \\ FB & 0.1 & & & & & 0.9 & \\ Sleep & & & & & & & 1 \end{array}$$
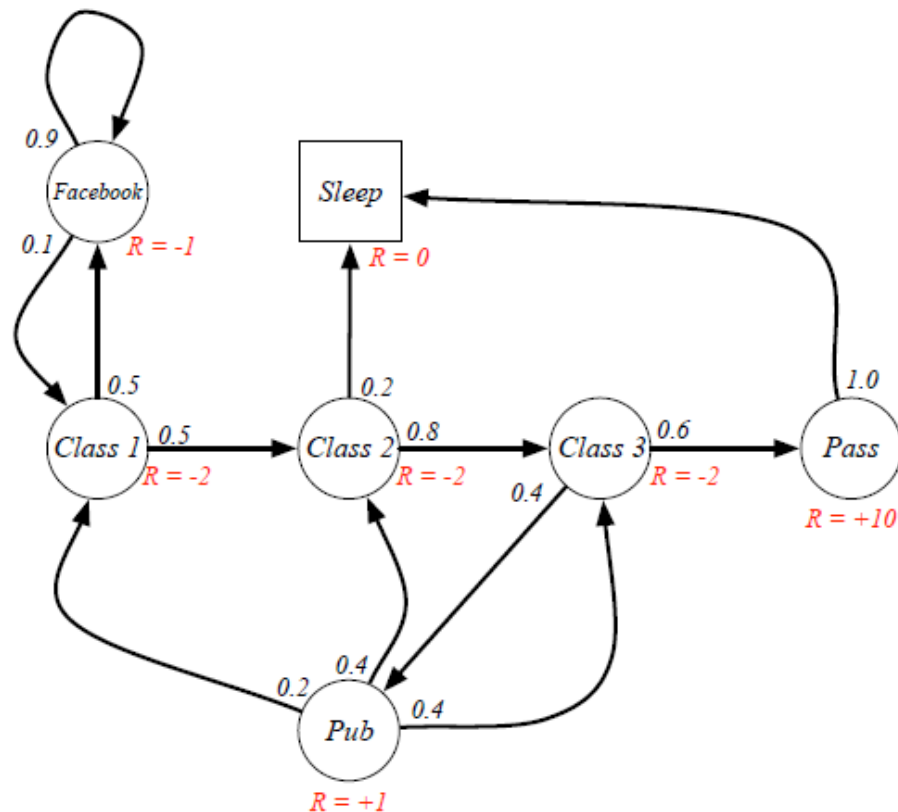
A Markov chain with **values**

**Markov Reward Process**

- A tuple $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
  - $\mathcal{S}$ is a (finite) set of states
  - $\mathcal{P}$ is a state transition probability matrix
  - $\mathcal{R}$ is a reward function, $\mathcal{R}_s = \mathbb{E}[R_{t+1}|S_t = s]$
  - $\gamma$ is a discount factor, $\gamma \in [0,1]$

Credits: David Silver's slides

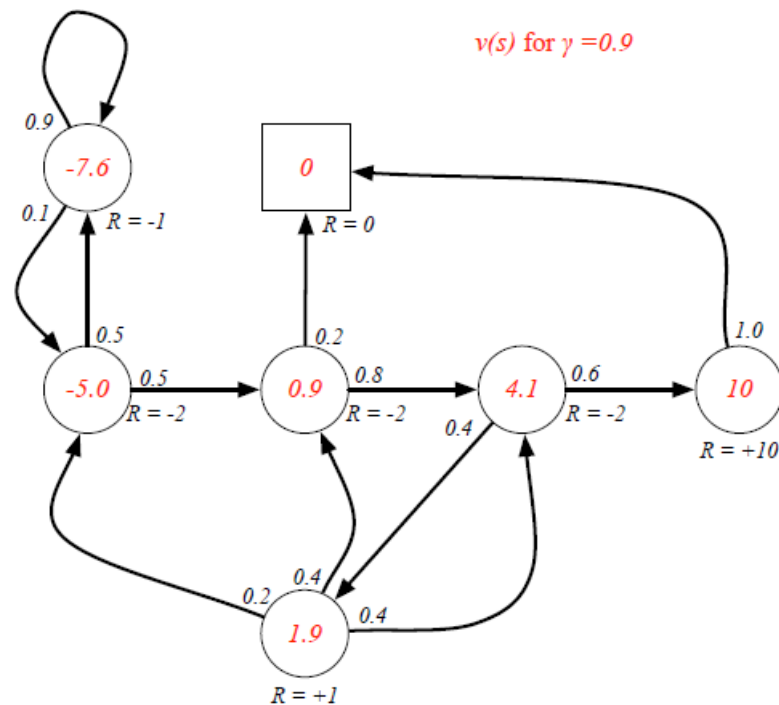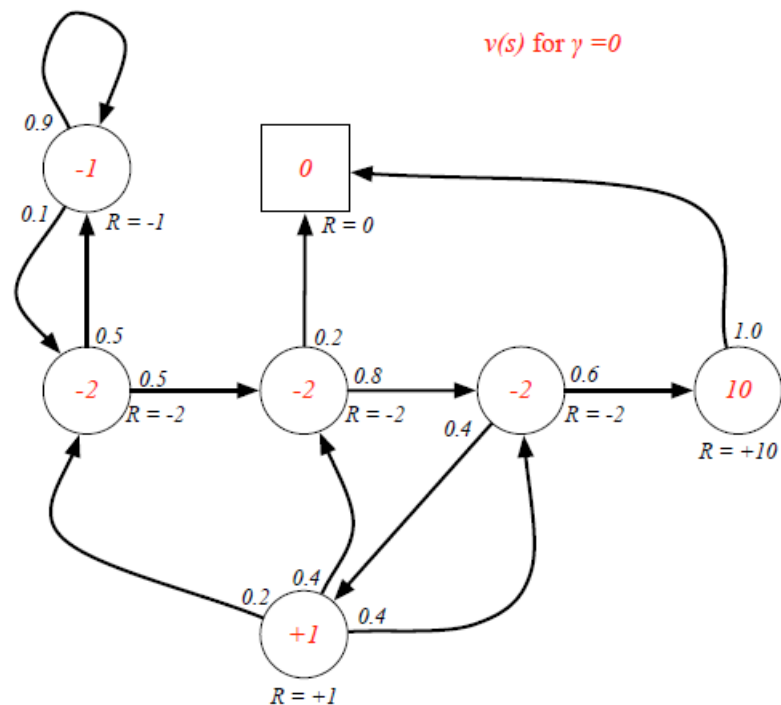Return $G_t$ is the **total discounted reward** from time-step $t$

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

➤ Discount $\gamma$ is the **present value of future rewards**
➤ Value immediate reward vs. delayed reward
  ➤ $\gamma$ near 0: "myopic"
  ➤ $\gamma$ near 1: "far-sighted"

The state value function $v(s)$ gives the expected return of state $s$

$$v(s) = \mathbb{E}[G_t | S_t = s]$$

Credits: David Silver's slides

The value function can be decomposed into **two parts**:
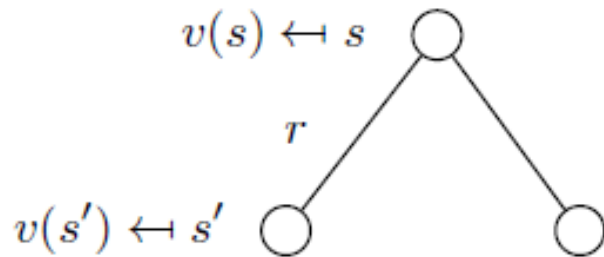- ▷ Immediate reward
- ▷ Discounted value of successor state

$$v(s) = \mathbb{E}[G_t | S_t = s]$$
$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots | S_t = s]$$
$$= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \gamma^2 R_{t+3} + \cdots) | S_t = s]$$
$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$
$$= \mathbb{E}[\textcolor{red}{R_{t+1}} + \textcolor{green}{\gamma v(S_{t+1})} | S_t = s]$$

**Immediate reward**

**Discounted value of next state**

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1})|S_t = s]$$



$v(s) \leftarrow s$

$r$

$v(s') \leftarrow s'$

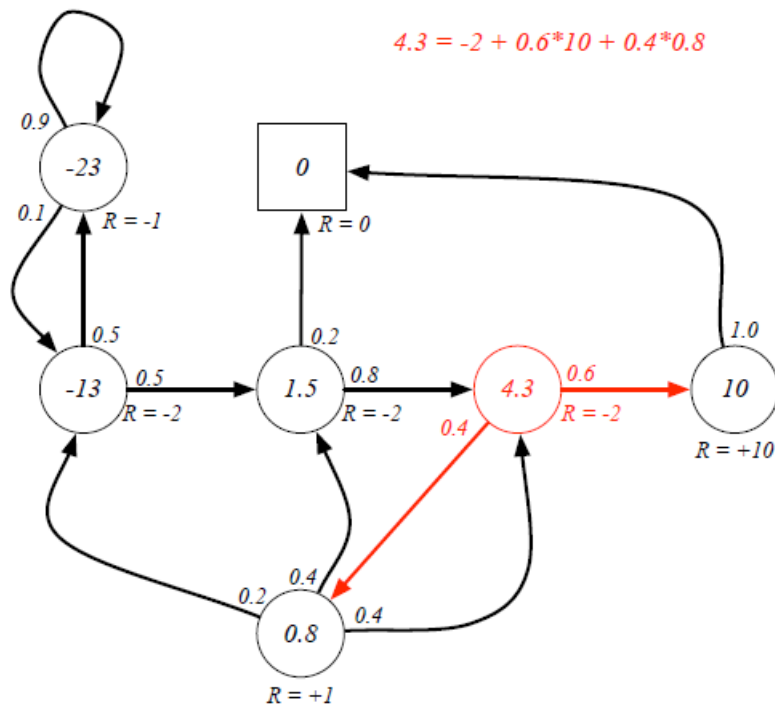$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} P_{ss'} v(s')$$

Credits: David Silver's slides

▷ $\mathcal{S}$ is a (finite) set of states
▷ $\mathcal{P}$ is a state transition probability matrix, with elements $P_{ss'}$
▷ $\mathcal{R}$ is a reward function, $\mathcal{R}_s = \mathbb{E}[R_{t+1}|S_t = s]$

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} \mathcal{R}(1) \\ \vdots \\ \mathcal{R}(n) \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{11} & \cdots & \mathcal{P}_{1n} \\ \vdots & \ddots & \vdots \\ \mathcal{P}_{n1} & \cdots & \mathcal{P}_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

$$v(s) = \mathcal{R} + \gamma \mathcal{P} v$$
$$(I - \gamma \mathcal{P}) v = \mathcal{R}$$
$$v = (I - \gamma \mathcal{P})^{-1} \mathcal{R}$$

- We can solve the Bellman equation! But it's $O(n^3)$ **for $n$ states**.
- For larger $n$ there are iterative solutions:
  - ➤ Monte-Carlo evaluation, Temporal-Difference learning
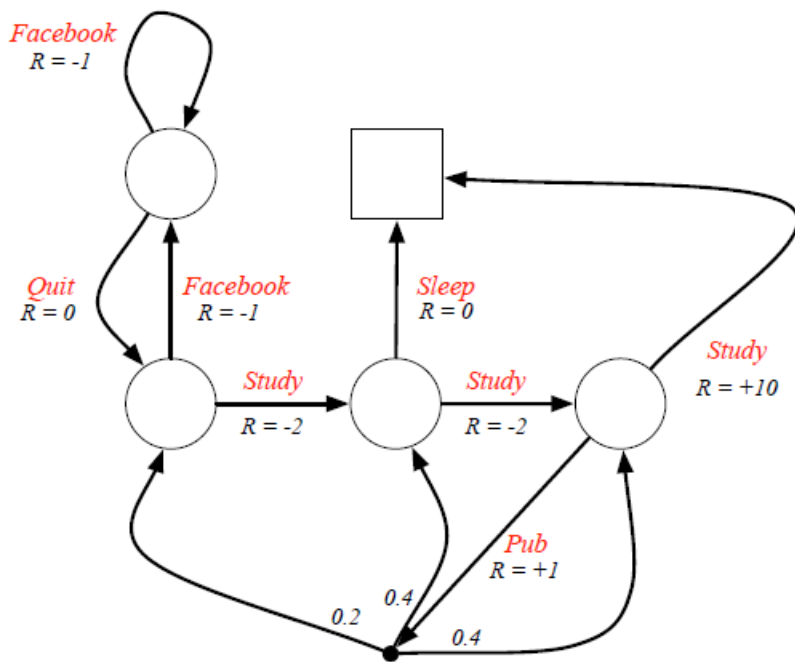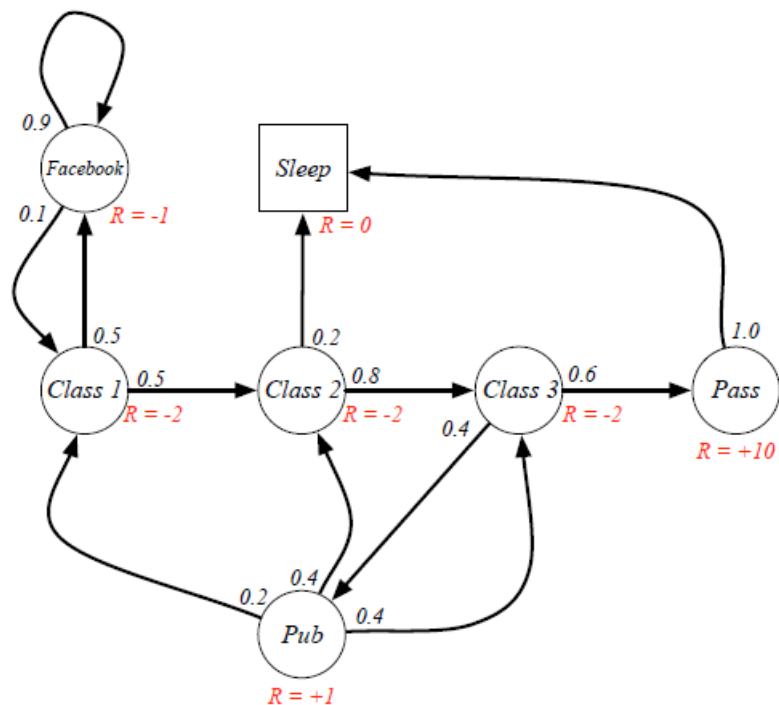  - ➤ Dynamic Programming

A Markov reward process **with decisions**.

**Markov Decision Process**

➤ A tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

➤ $\mathcal{S}$ is a finite set of states

➤ $\mathcal{A}$ is a finite set of actions

➤ $\mathcal{P}$ is a state transition probability matrix
$$\mathcal{P}_{ss'}^{a} = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

➤ $\mathcal{R}$ is a reward function, $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$

➤ $\gamma$ is a discount factor, $\gamma \in [0,1]$

Credits: David Silver's slides

**Policy** ($\pi$): a distribution of actions given states
- $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$
- Fully defines the behavior of an agent
- MDP policies depend on the current state, not the history
- *Stationary* (time-independent), since you're considering all possible states

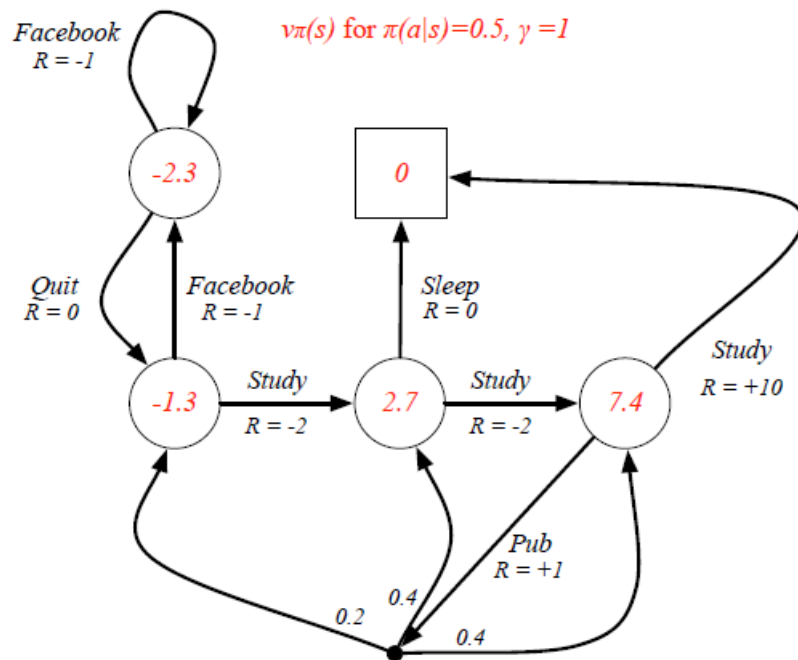State-value function $(v_\pi(s))$: the expected return following policy $\pi$, from state $s$

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

Action-value function $(q_\pi(s, a))$: the expected return after **taking action $a$**, following policy $\pi$, from state $s$

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

Credits: David Silver's slides

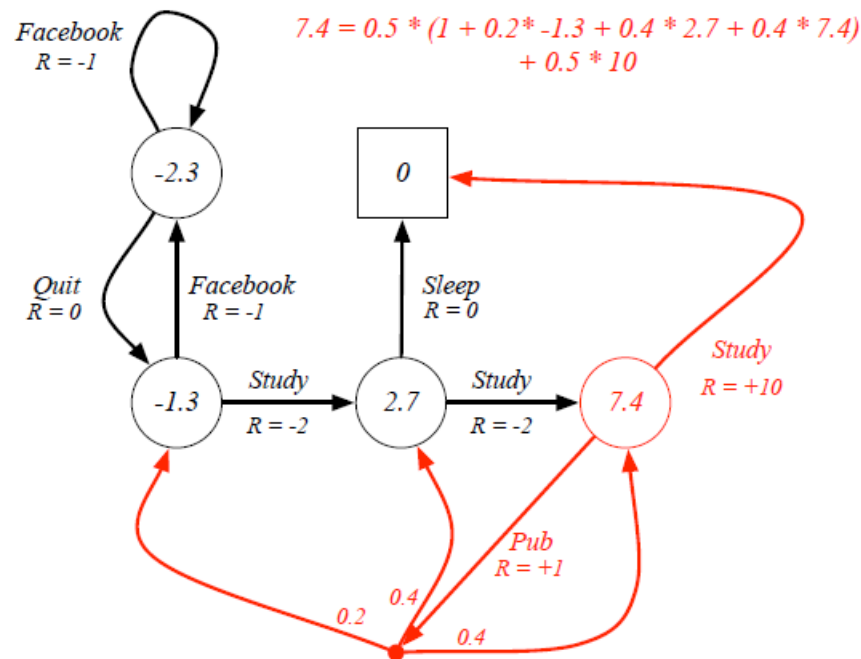**State-value function** $(v_\pi(s))$: the expected return following policy $\pi$, from state $s$

$$v_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1})|S_t = s]$$

**Action-value function** $(q_\pi(s, a))$: the expected return after **taking action $a$**, following policy $\pi$, from state $s$

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

$7.4 = 0.5 * (1 + 0.2 * -1.3 + 0.4 * 2.7 + 0.4 * 7.4)$
$+ 0.5 * 10$

Facebook
R = -1

-2.3

0

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

-1.3   Study   2.7   Study   7.4
       R = -2         R = -2

Pub
R = +1

0.4

0.2

0.4

Credits: David Silver's slides

■ **Optimal state-value function** $(v_*(s))$: the maximum state-value function over all policies

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

■ **Optimal action-value function** $(q_*(s, a))$: the maximum action-value function over all policies

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

■ Theorem:
  ➤ There exists an optimal policy $\pi_*$ better than or equal to all other policies, $\pi_* \geq \pi, \forall \pi$
  ➤ All optimal policies achieve the optimal state-value and action-value functions
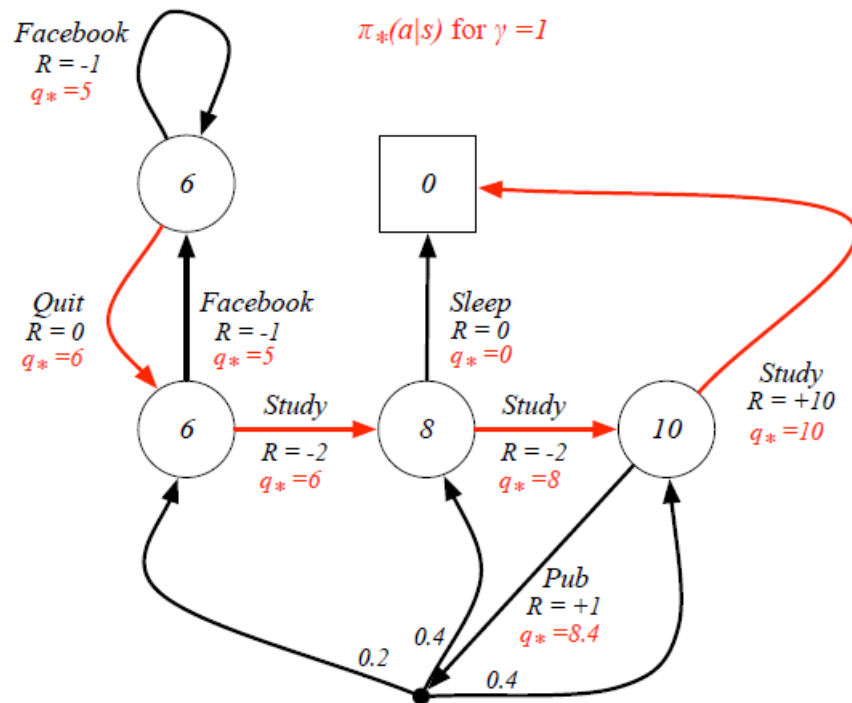
An optimal policy can be found by maximizing over $q_*(s, a)$,

$$\pi_*(a|s) = \begin{cases} 1, & \text{if } a = \operatorname*{argmax}_{a \in \mathcal{A}} q_*(s, a) \\ 0, & \text{otherwise} \end{cases}$$

Credits: David Silver's slides

The Bellman equation for the state and action-value functions

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

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

# Solving the Bellman Optimality Equation

- Bellman Optimality Equation is non-linear
- No closed form solution (in general)
- Iterative solution methods:
  - Value Iteration
  - Policy Iteration
  - Q-learning
  - Sarsa

# Thank you!

# Sources

- [UCL Course on Reinforcement Leaning](), David Silver

- [Reinforcement Learning – An Introduction](), Sutton and Barto

# ♡ CREDITS

Special thanks to:

- Presentation template by SlidesCarnival
- Photographs by Startup Stock Photos