

Lecture 1: Introduction to RL

Zhi Wang

Department of Control and Systems Engineering
Nanjing University

Nov. 5th, 2019

Table of Contents

- 1 Overview of RL
- 2 Finite Markov Decision Process

What we'll cover

- From supervised learning to decision making
- Tabular RL algorithms:
 - Dynamic programming
 - Monte-Carlo methods
 - Temporal difference learning
- Deep RL algorithms:
 - Policy gradients
 - Actor-critics
 - Value function-based methods
- State-of-the-arts:
 - Transfer RL, meta-RL
 - Evolution strategies for RL
 - Model-based RL
 - Hierarchical RL
 - Inverse RL

Useful links and references

- Course website
 - http://heyuanmingong.github.io/rl_lecture/index.html
- Book
 - Richard S. Sutton, and Andrew G. Barto, *Reinforcement Learning: An Introduction*. 2nd Edition.
- Famous open classes
 - CS 285 at UC Berkeley, Deep Reinforcement Learning,
<http://rail.eecs.berkeley.edu/deeprlcourse/>
 - CS 234 at Stanford, Reinforcement Learning,
<http://web.stanford.edu/class/cs234/>

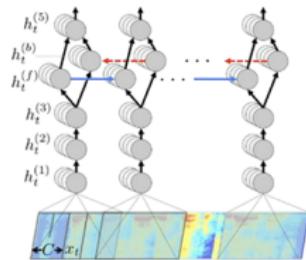
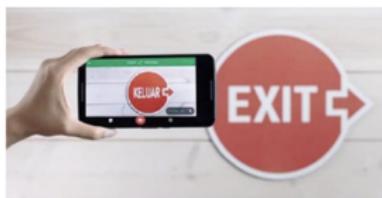
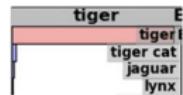
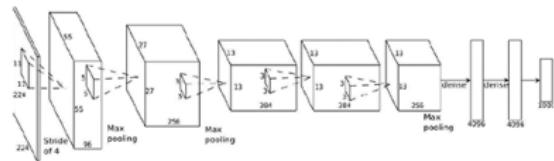
What is RL, and why we should care?

- Intelligent machines must able to adapt



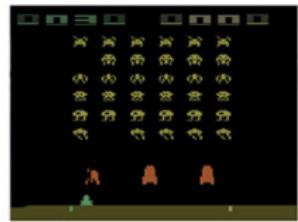
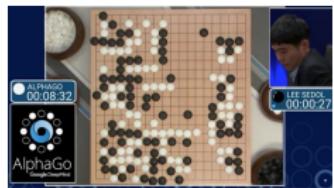
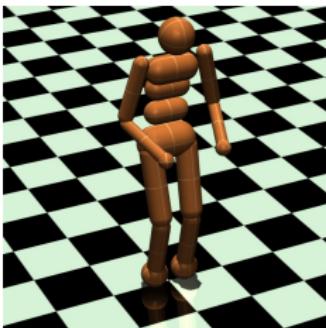
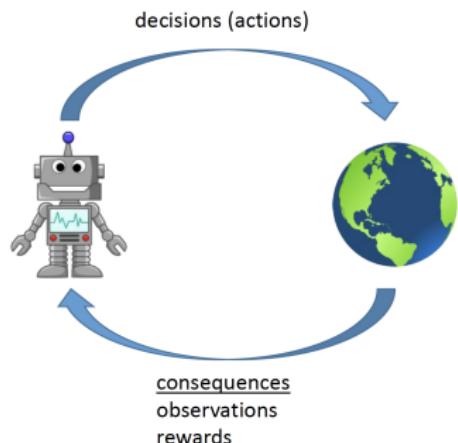
What is RL, and why we should care?

- Deep learning helps us handle *unstructured environments*



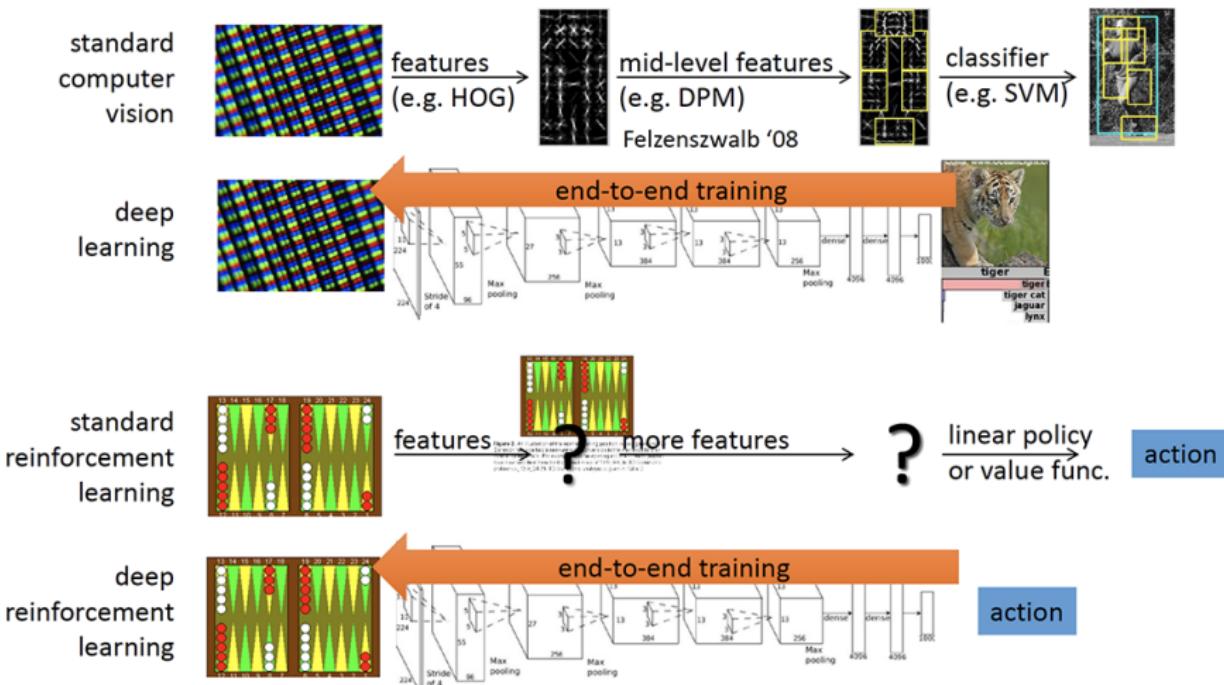
What is RL, and why we should care?

- RL provides a formalism for behavior, decision-making



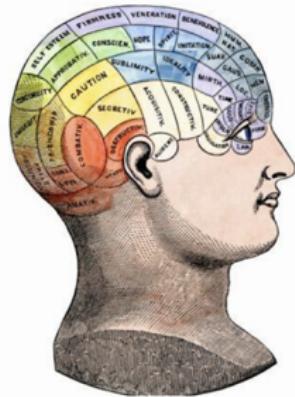
Space Invaders

How to use RL?

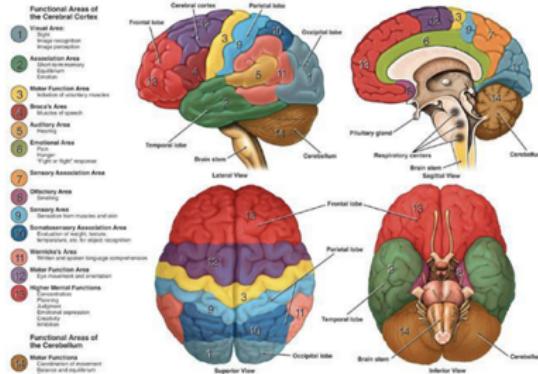


How do we build intelligent machines?

- Imagine you have to build an intelligent machine, where do you start?



Anatomy and Functional Areas of the Brain

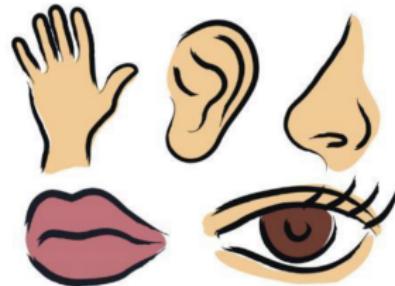


Learning as the basis of intelligence

- Some things we can all do (e.g. walking)
- Some things we can only learn (e.g. driving a car)
- We can learn a huge variety of things, including very difficult things
- Therefore our learning mechanism(s) are likely powerful enough to do everything we associate with intelligence
 - But it may still be very convenient to “hard code” a few really important bits

What must a single algorithm do?

- Interpret (rich) sensory inputs



- Choose (complex) actions



Reinforcement learning in the brain

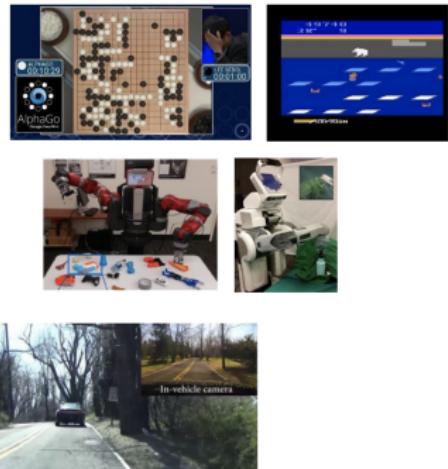
Yael Niv

Psychology Department & Princeton Neuroscience Institute, Princeton University

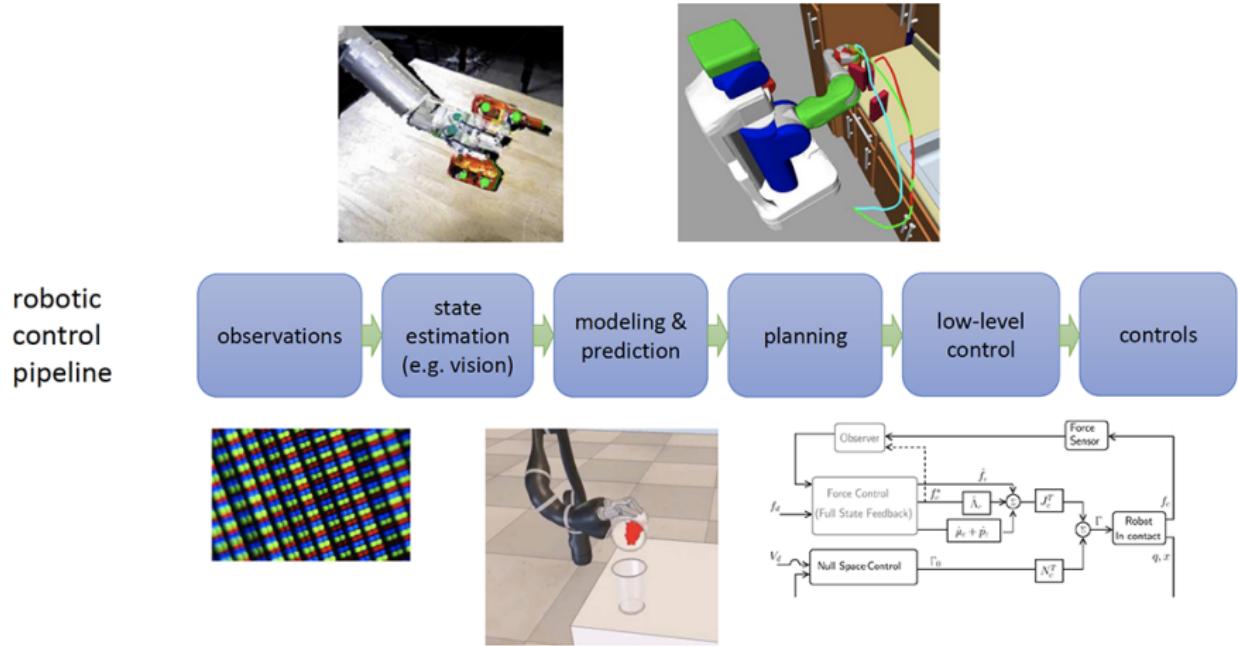
- Percepts that anticipate reward become associated with similar firing patterns as the reward itself
- Basal ganglia appears to be related to reward system
- Model free RL like adaptation is often a good fit for experimental data of animal adaptation (but not always)

What can (deep) RL do?

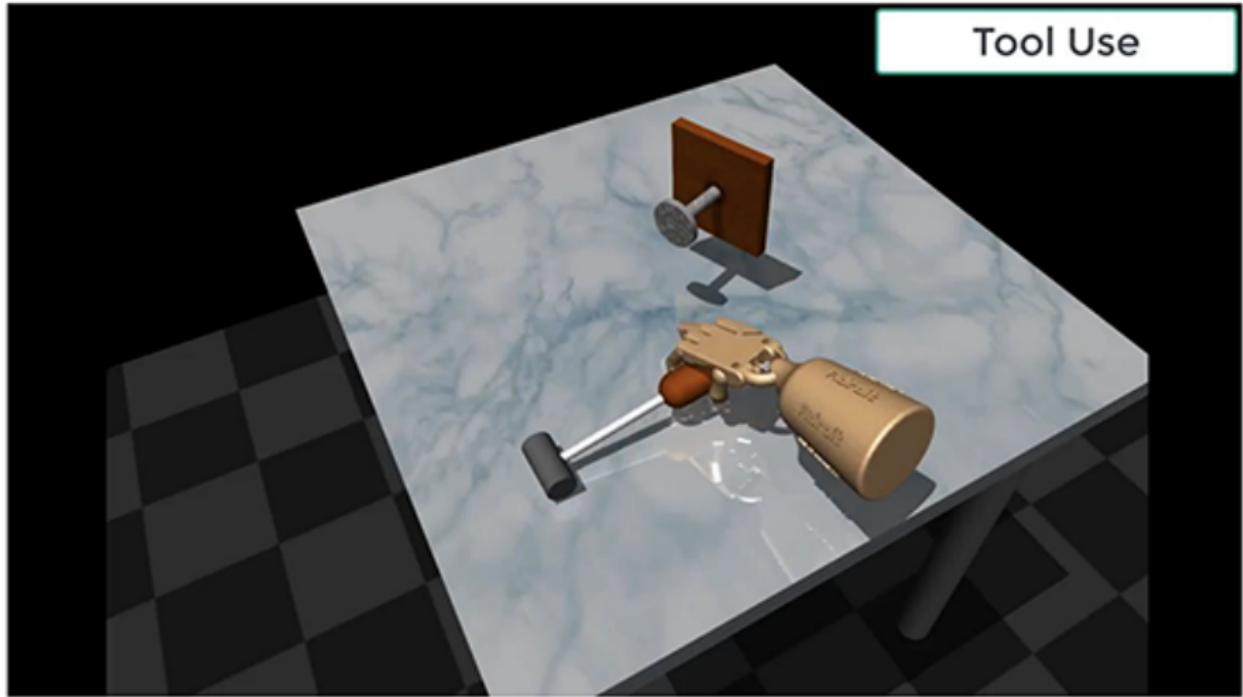
- Acquire high degree of proficiency in domains governed by simple, known rules
- Learn simple skills with raw sensory inputs, given enough experience
- Learn from imitating enough human provided expert behavior



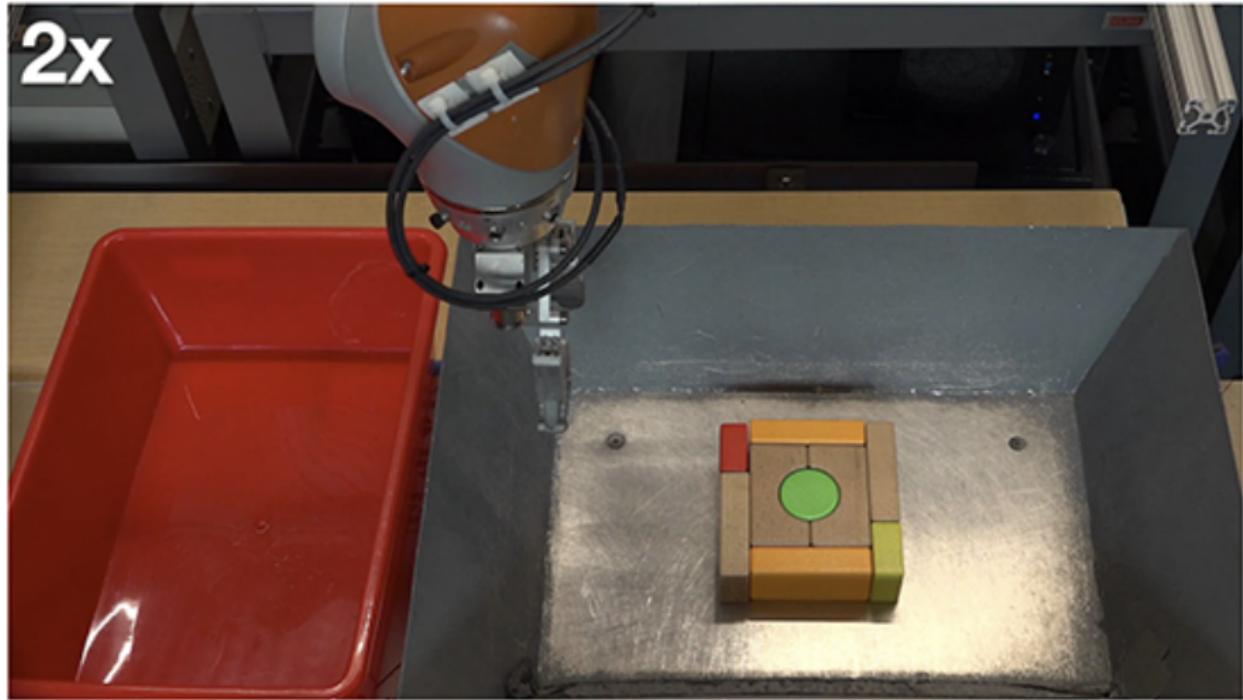
Example: Robotics



Example: Complex physical tasks



Example: Real-world scenarios



Why should we study this now?



Atari games:

Q-learning:

V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, et al. "Playing Atari with Deep Reinforcement Learning". (2013).

Policy gradients:

J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. "Trust Region Policy Optimization". (2015).
V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, et al. "Asynchronous methods for deep reinforcement learning". (2016).

Real-world robots:

Guided policy search:

S. Levine*, C. Finn*, T. Darrell, P. Abbeel. "End-to-end training of deep visuomotor policies". (2015).

Q-learning:

D. Kalashnikov et al. "QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation". (2018).

Beating Go champions:

Supervised learning + policy gradients + value functions + Monte Carlo tree search:

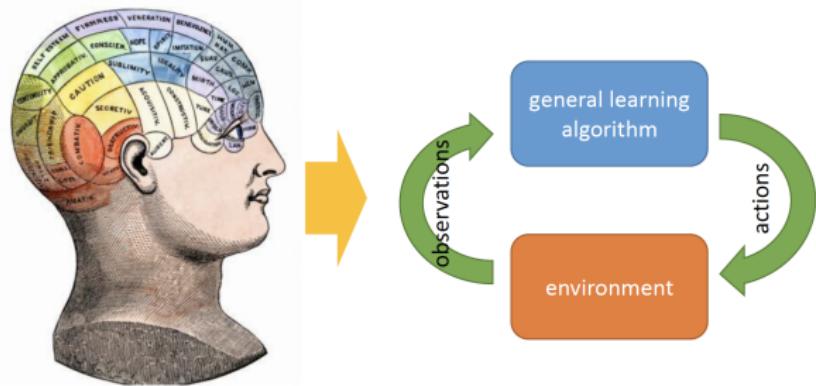
D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, et al. "Mastering the game of Go with deep neural networks and tree search". Nature (2016).

RL, decision making, human behavior

Instead of trying to produce a program to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain.

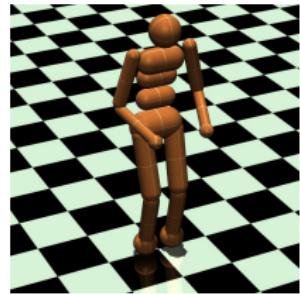
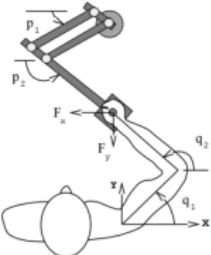


- Alan Turing



What has proven challenging so far?

- Humans can learn incredibly quickly
 - (Deep) RL methods are usually slow
- Humans can reuse past knowledge
 - Transfer learning in (deep) RL is an open problem
- Not clear what the reward function should be
- Not clear what the role of prediction should be



To enable real-world sequential decision making – Beyond learning from reward

- Basic RL deals with maximizing rewards
- This is not the only problem that matters for sequential decision making!
- More advanced topics
 - Transferring knowledge between domains (transfer learning, meta learning)
 - Learning to predict and using prediction to act (model-based RL)
 - Learning reward functions from example (inverse RL)

Table of Contents

- 1 Overview of RL
- 2 Finite Markov Decision Process

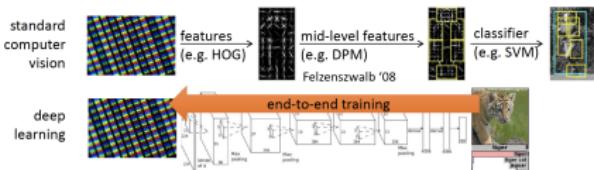
Markov Decision Process (MDP)

- A classical formalization of sequential decision making
 - Choosing different actions in different situations
 - Actions influence not just immediate rewards, but also subsequent situations through future rewards
 - Involve delayed reward and the need to tradeoff immediate and delayed reward
- A mathematically idealized form of the RL problem
 - Precise theoretical statements can be made
 - Key elements of the problem's mathematical structure, such as returns, value functions, Bellman equations, etc
 - A tension between breadth of applicability and mathematical tractability

Supervised learning vs. Sequential decision making

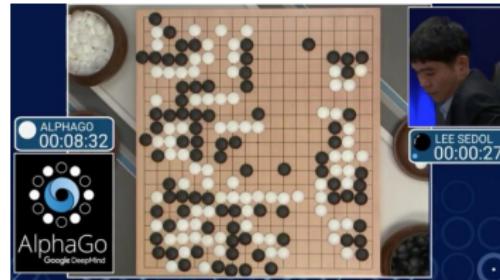
Supervised learning

- Samples are independent and identically distributed (i.i.d.)
- Given an input, map an optimal output

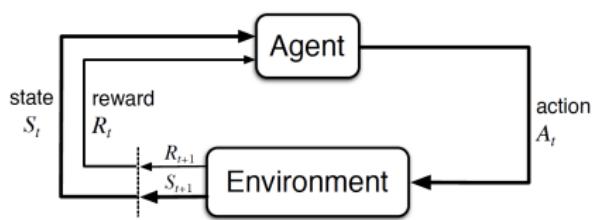


Reinforcement learning

- Samples are not i.i.d., temporally co-related
- Given an initial state, find a sequence of optimal actions



The agent-environment interface



- **Agent:** The learner, decision maker
- **Environment:** The thing it interact with, comprising everything outside the agent

At each time step $t = 0, 1, 2, \dots$, the agent...

- receives some representation of the environment's **state**, $S_t \in \mathcal{S}$
- on that basis, selects an **action**, $A_t \in \mathcal{A}(s)$
- one time step later, receives a numerical **reward**, $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$
- finds itself in a new state, S_{t+1} , (**transition function**)
- $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3 \dots$

Dynamics of the MDP $p : \mathcal{S} \times \mathcal{R} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$

$$p(s', r|s, a) = \Pr\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r|s, a) = 1, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$$

- The probabilities given by p completely characterize the environment's dynamics
- **Markov Property**
 - The probability of each possible value for S_t and R_t depends only on the immediately preceding state S_{t-1} and action A_{t-1} , not at all on earlier states and actions
 - $p(S_t, R_t | S_{t-1}, A_{t-1}) = p(S_t, R_t | S_{t-1}, A_{t-1}, S_{t-2}, A_{t-2}, S_{t-3}, A_{t-3}, \dots)$
- Recall supervised learning $p(X_i | X_j) = 0$

- **Reward Hypothesis**

- That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).
- The agent's goal: maximize the total amount of reward it receives
 - $\maximize J(\pi) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r_{t+1}]$
 - Maximize not immediate reward, but cumulative reward in the long run

Returns and episodes

- The subsequence of the agent-environment interaction, **episodes**
 - $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
 - Each episode ends in a special state called the terminal state, S_T
- We seek to maximize the **expected return**, G_t , the reward sequence
 - $G_t = R_{t+1} + R_{t+2} + \dots + R_T$
 - Episodic tasks

Discount rate $\gamma \in [0, 1]$

- The expected **discounted** return
 - $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots = \sum_{k=1}^{\infty} \gamma^k R_{t+k+1}$
- The discount rate determines the present value of future rewards: a reward received k time steps in the future is worth only γ^{k-1} times what it would be worth if it were received immediately
- $\gamma \rightarrow 0$, the agent is “myopic”, only maximizing immediate rewards
 - Akin to supervised learning that maximizes the log-likelihood of each sample, $\log p(y_i|x_i)$
- $\gamma \rightarrow 1$, the agent is “farsighted”, taking future rewards into account
- Returns at successive time steps are related to each other

$$\begin{aligned} G_t &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

Policies and value functions

- **Value functions:** function of states or state-action pairs
 - Estimate how good it is for the agent to be in a given state, or how good it is to perform a given action in a given state
 - “How good”: defined in expected future rewards, i.e., expected return
 - Depend on what actions to take, defined w.r.t. particular ways of acting, called **policies**
- **Policy:** $\pi(a|s)$, a mapping from states to probabilities of selecting each possible action
 - $\sum_a \pi(a|s) = 1$
 - e.g., for a given state s_1 , four possible actions
 $a_1 = 0.1, a_2 = 0.3, a_3 = 0.2, a_4 = 4$
- RL = estimate value functions + estimate policies + estimate both

Policies and value functions

- v_π , the **state-value function** for policy π

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi \left[\sum_{k=1}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right], \forall s \in \mathcal{S}$$

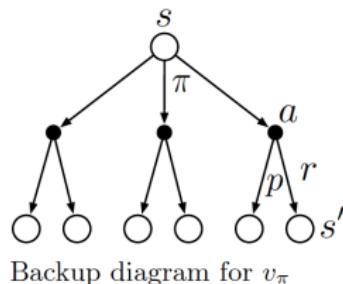
- Q_π , the **action-value function** for policy π

$$\begin{aligned} Q_\pi(s, a) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi \left[\sum_{k=1}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right], \forall s \in \mathcal{S}, a \in \mathcal{A}(s) \end{aligned}$$

A fundamental property: Bellman equation

$$\begin{aligned}v_{\pi}(s) &= \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s']] \\&= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]\end{aligned}$$

- Average over all the possibilities, weighting each by its probability of occurring
- Express the relationship between the value of a state and the values of its successor states
- Transfer value information back to a state from its successor states



Optimal policies and optimal value functions

- RL tasks: find a policy that achieves a lot of reward over the long run
- Value functions define a partial ordering over policies
 - $\pi \geq \pi'$ if and only if $v_\pi(s) \geq v_{\pi'}(s), \forall s \in \mathcal{S}$
- At least one policy that is better than or equal to all other policies, i.e., **optimal policy**
- Optimal policies, π^* , share the same **optimal value function**

$$v^*(s) = \max_{\pi} v_{\pi}(s), \quad \forall s \in \mathcal{S}$$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$$

$$Q^*(s, a) = \mathbb{E}[R_{t+1} + \gamma v^*(S_{t+1}) | S_t = s, A_t = a]$$

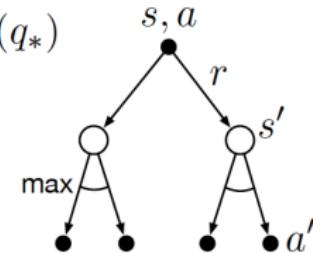
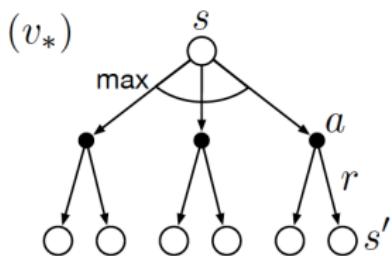
Bellman optimality equation

$$\begin{aligned}v^*(s) &= \max_a Q_{\pi^*}(s, a) \\&= \max_a \mathbb{E}_{\pi^*}[G_t | S_t = s, A_t = a] \\&= \max_a \mathbb{E}_{\pi^*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= \max_a \mathbb{E}[R_{t+1} + \gamma v^*(S_{t+1}) | S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v^*(s')]\end{aligned}$$

- The value of a state under an optimal policy must equal the expected return for the best action from that state

Bellman optimality equation

$$\begin{aligned} Q^*(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_{a'} Q^*(S_{t+1}, a') | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q^*(s', a')] \end{aligned}$$



$$v^*(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v^*(s')]$$

Determine optimal policies from optimal value functions

- For v^* : a one-step search
 - Actions that appear best after one-step search will be optimal actions
- For Q^* : no need to do a one-step-ahead search
 - $a^* = \arg \max_a Q^*(s, a)$
 - The optimal action-value function allows optimal actions to be selected without having to know anything about possible successor states and their values, i.e., without having to know anything about the environment's dynamics