

CS 291A: Deep Learning for NLP

Deep Reinforcement Learning (1)

William Wang
UCSB Computer Science
william@cs.ucsb.edu

Slides adapted from V. Chen and D. Silver.

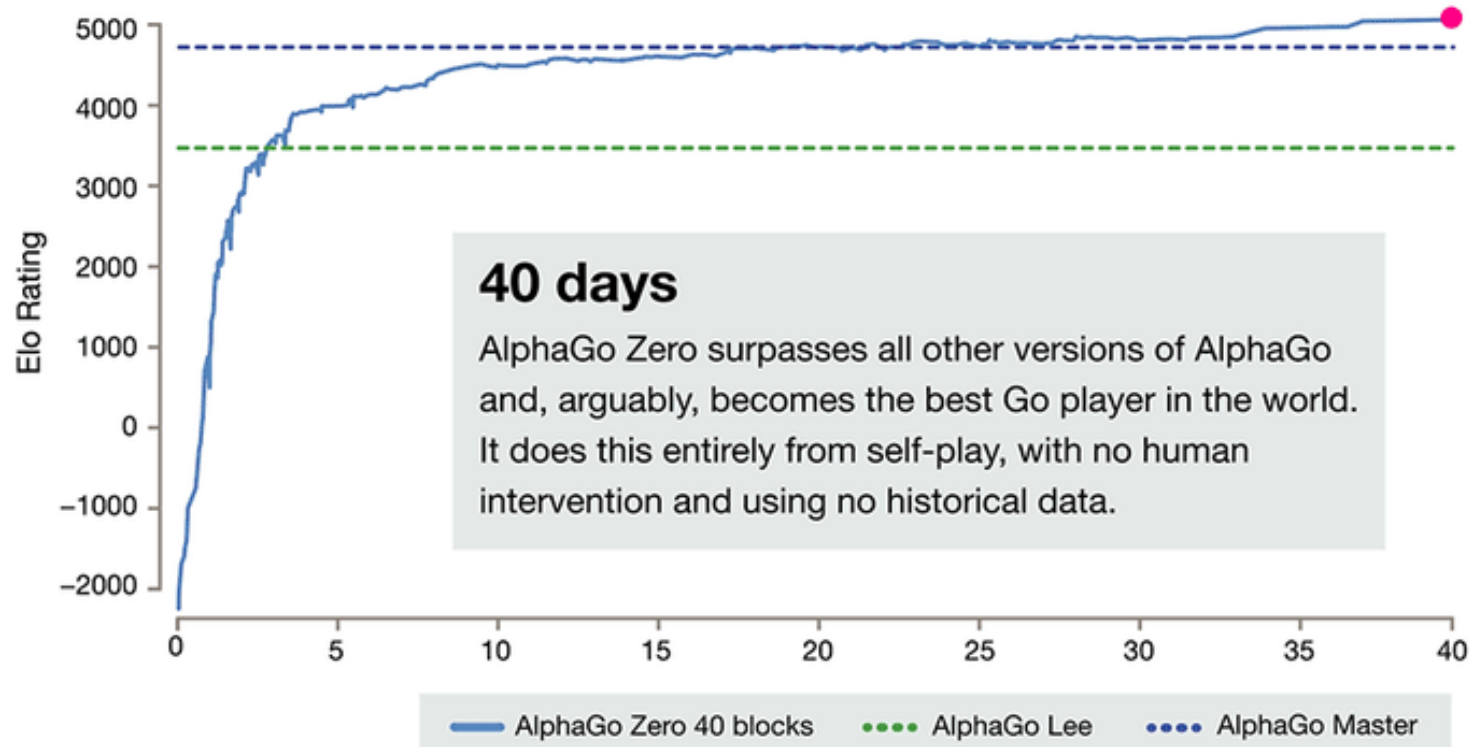
Course Progress

- We are done with homework assignments.
- We are less than two weeks from the scheduled final project presentations.
 - Three batches of final project presentation (schedule will be announced shortly).
 - We expect results, not just ideas.
 - If you need GPU instances on Google Cloud, you need to request the limit increase through their customer service (same as AWS).

New **AlphaGO** (based on RL) beats World's No. 1 Go player **Ke Jie** last year.



AlphaGo Zero



Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

Reinforcement Learning Approach

- Policy-Based
- Value-Based
- Model-Based

Problems within RL

- Learning and Planning
- Exploration and Exploitation

Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

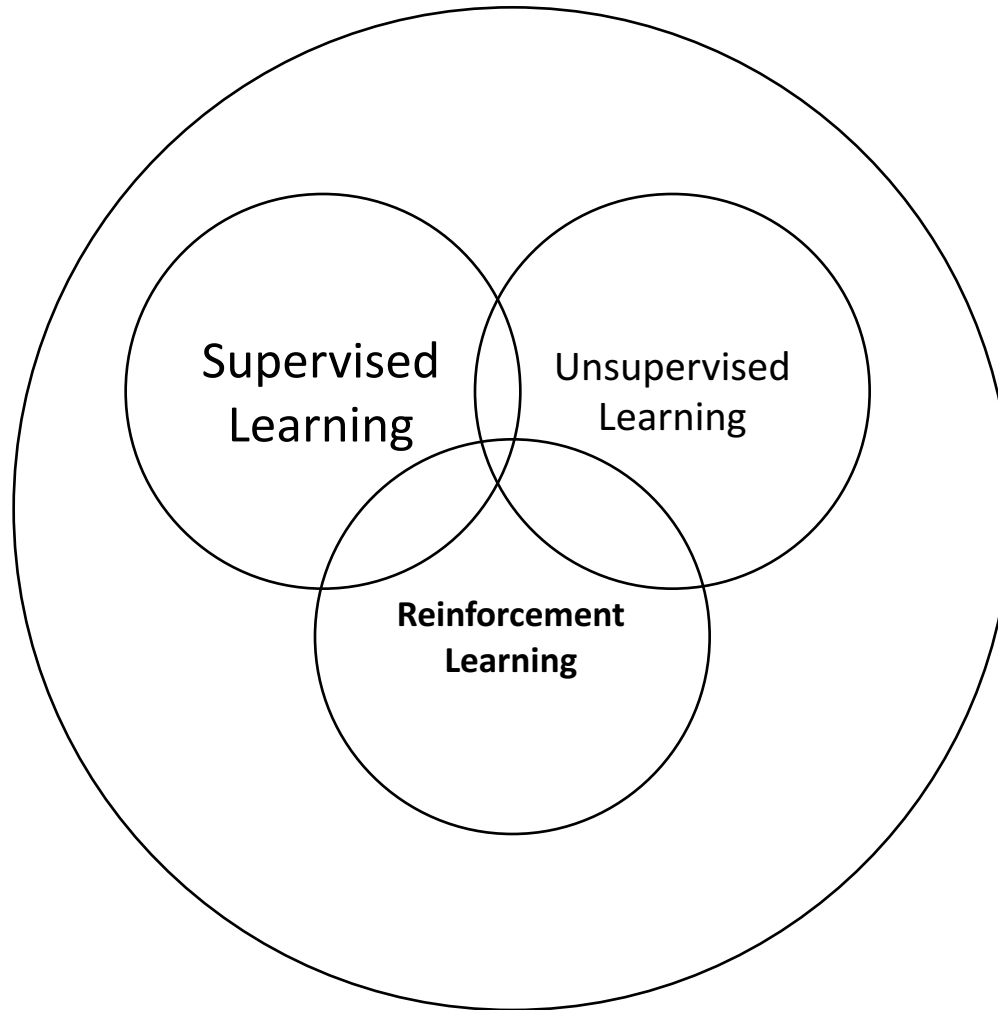
Reinforcement Learning Approach

- Policy-Based
- Value-Based
- Model-Based

Problems within RL

- Learning and Planning
- Exploration and Exploitation

Machine Learning



Supervised v.s. Reinforcement

Supervised Learning

- Training based on supervisor/label/annotation
- Feedback is instantaneous
- Not much temporal aspects

Reinforcement Learning

- Training only based on reward signal
- Feedback is delayed
- Time matters
- Agent actions affect subsequent exploration

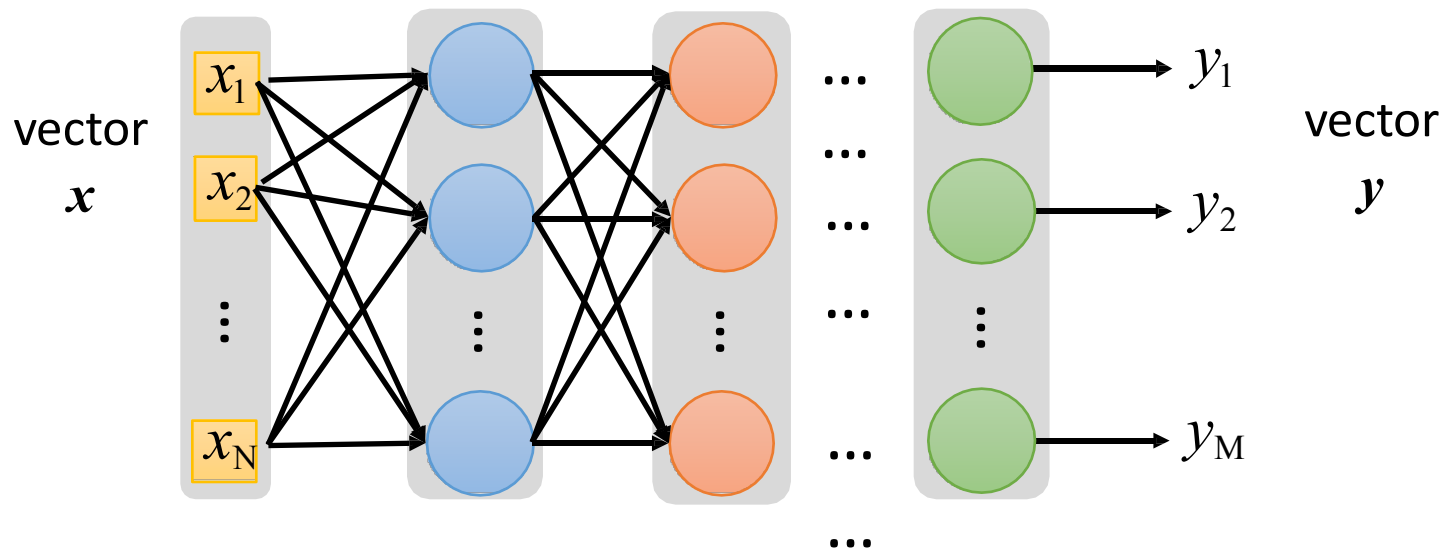
Reinforcement Learning

- RL is a general purpose framework for **decision making**
- ◦ RL is for an *agent* with the capacity to *act*
- ◦ Each *action* influences the agent's future *state*
- ◦ Success is measured by a scalar *reward* signal
- ◦ Goal: *select actions to maximize future reward*

Deep Learning

DL is a general purpose framework for **representation learning**

- Given an *objective*
- Learn *representation* that is required to achieve objective
- Directly from *raw inputs*
- Use minimal domain knowledge



Deep Reinforcement Learning

- DRL is an agent that can solve human-level task
 - ◦ RL defines the objective
 - ◦ DL learns representations
 - ◦ RL + DL = artificial intelligence

Deep RL AI Examples

- Play games: Atari, poker, Go, ... Explore worlds: 3D worlds, ...
- Control physical systems: manipulate, ...
- Interact with users: recommend, optimize, personalize, ...



Introduction to RL

Reinforcement Learning

Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

Reinforcement Learning Approach

- Policy-Based
- Value-Based
- Model-Based

Problems within RL

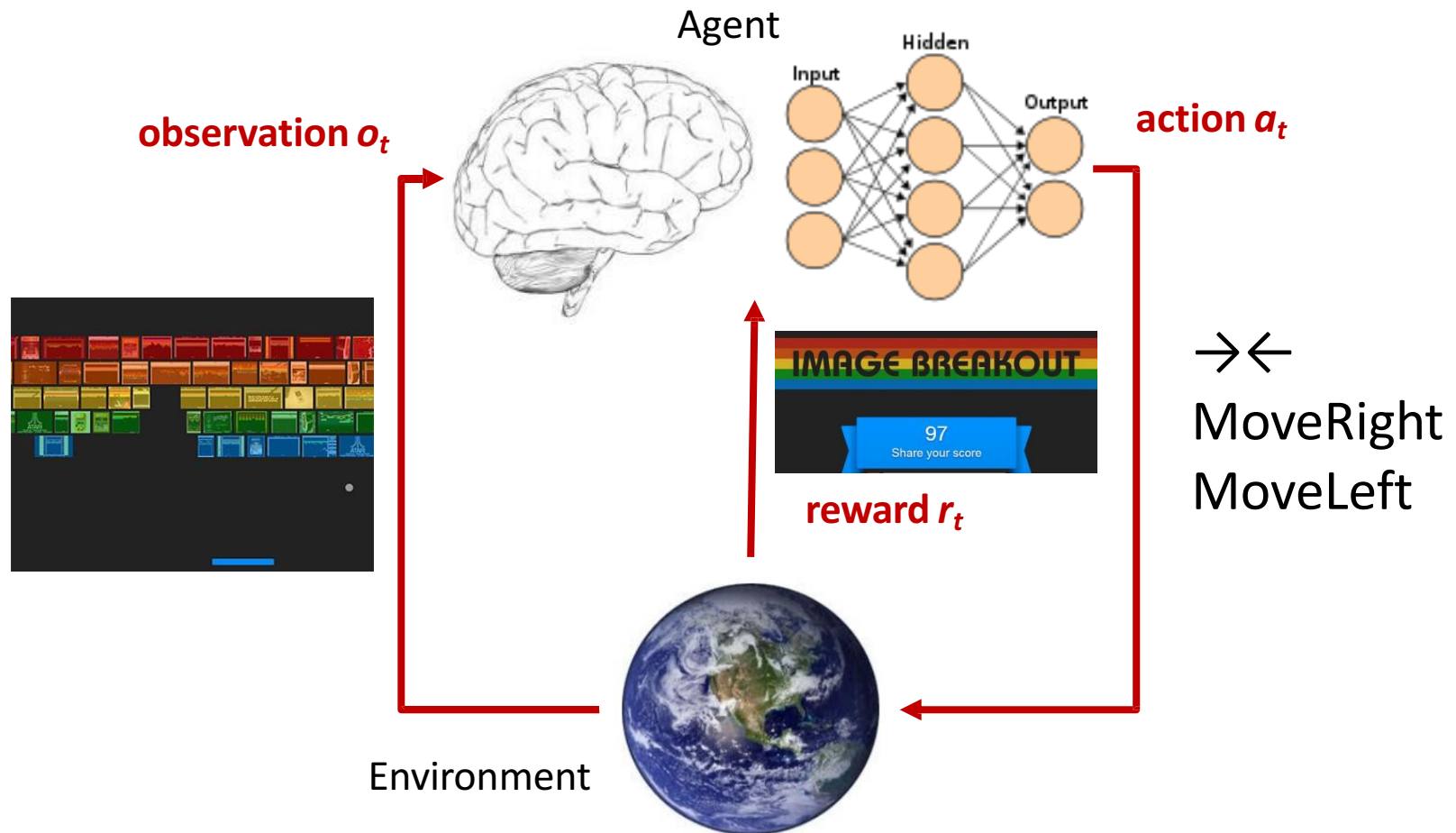
- Learning and Planning
- Exploration and Exploitation

Reinforcement Learning

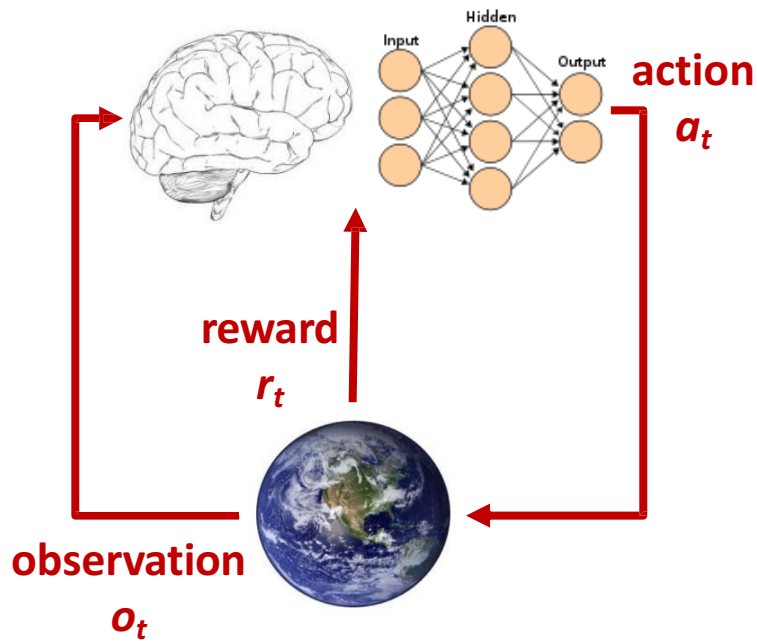
- RL is a general purpose framework for **decision making**
- ◦ RL is for an *agent* with the capacity to *act*
- ◦ Each *action* influences the agent's future *state*
- ◦ Success is measured by a scalar *reward* signal

Big three: action, state, reward

Agent and Environment



Agent and Environment



At time step t

- The agent
 - Executes action a_t
 - Receives observation o_t
 - Receives scalar reward r_t
- The environment
 - Receives action a_t
 - Emits observation o_{t+1}
 - Emits scalar reward r_{t+1}
- t increments at env.step

State

Experience is the sequence of observations, actions, rewards

$$o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t$$

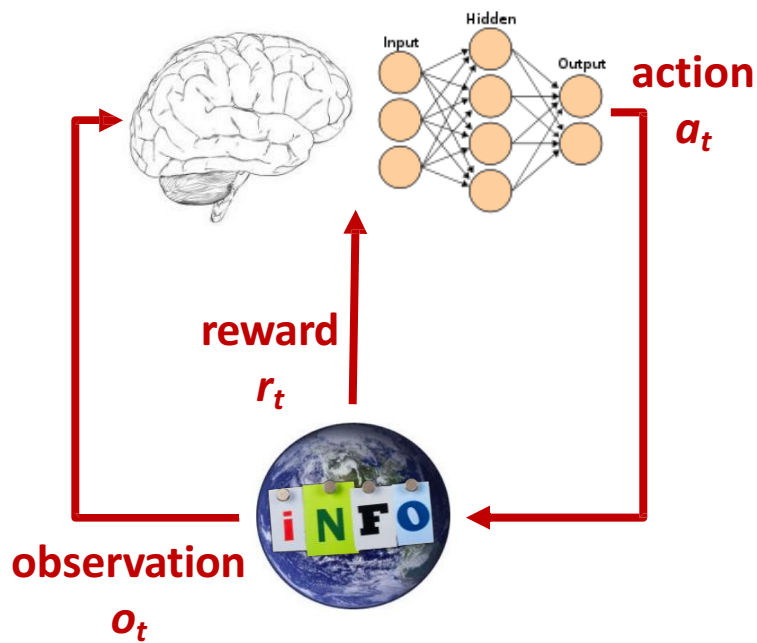
State is the information used to determine what happens next

- what happens depends on the history experience
 - The agent selects actions
 - The environment selects observations/rewards

The state is the function of the history experience

$$s_t = f(o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t)$$

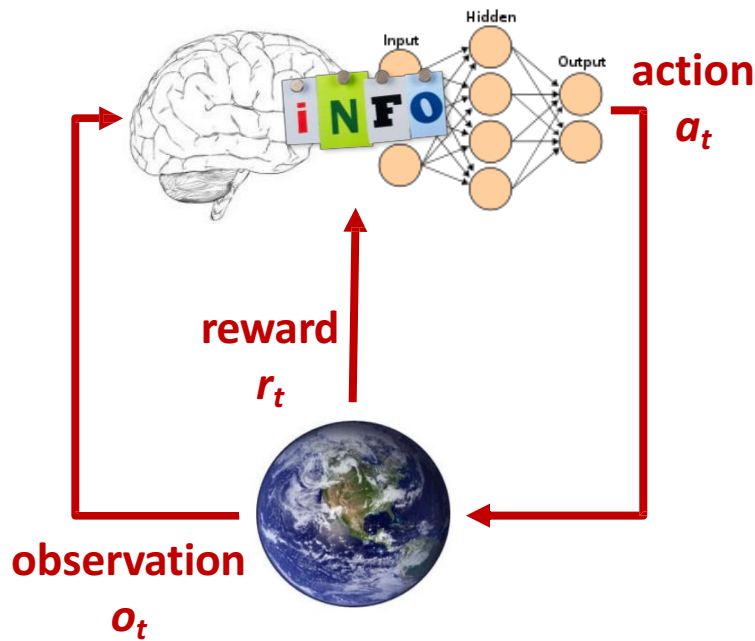
Environment State



The **environment state** s_t^e is the environment's *private* representation

- impact on the next observation / reward
- may not be visible to the agent
- may contain irrelevant information

Agent State



The **agent state** s_t^a is the agent's *internal* representation

- whether data the agent uses to pick the next action \rightarrow information used by RL algorithms
- can be any function of experience

Information State

An information state (a.k.a. Markov state) contains all useful information from history

A state is Markov iff $P(s_{t+1} \mid s_t) = P(s_{t+1} \mid s_1, \dots, s_t)$

The future is independent of the past given the present

$$H_t = \{o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t\}$$

$$H_{1:t} \rightarrow s_t \rightarrow H_{t+1:\infty}$$

- Once the state is known, the history may be thrown away
- The state is a sufficient statistics of the future

Fully Observable Environment

Full observability: agent directly observes environment state

$$O_t = s_t^a = s_t^e$$

information state = agent state = environment state

This is a Markov decision process (MDP)

Partially Observable Environment

Partial observability: agent indirectly observes environment

$$s_t^a \neq s_t^e$$

agent state \neq environment state

This is partially observable Markov decision process (POMDP)

Agent must construct its own state representation s_t^a

- Complete history: $s_t^a = H_t$
- Beliefs of environment state: $s_t^a = \{P(s_t^e = s^1), \dots, P(s_t^e = s^n)\}$
- Hidden state (from RNN): $s_t^a = \sigma(W_s \cdot s_{t-1}^a + W_o \cdot o_t)$

Reward

Reinforcement learning is based on reward hypothesis

A reward r_t is a scalar feedback signal

- Indicates how well agent is doing at step t

Reward hypothesis: all agent goals can be desired by maximizing expected cumulative reward

Sequential Decision Making

Goal: select actions to maximize total future reward

- Actions may have long-term consequences
- Reward may be delayed
- It may be better to sacrifice immediate reward to gain more long-term reward



Markov Decision Process

Fully Observable Environment

Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

Reinforcement Learning Approach

- Policy-Based
- Value-Based
- Model-Based

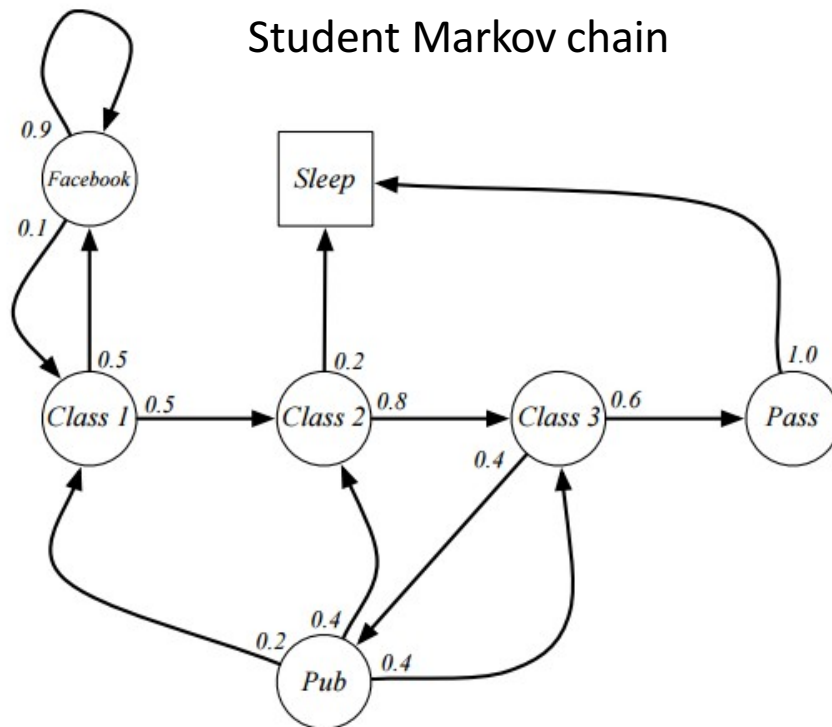
Problems within RL

- Learning and Planning
- Exploration and Exploitation

Markov Process

Markov process is a memoryless random process

- i.e. a sequence of random states S_1, S_2, \dots with the Markov property



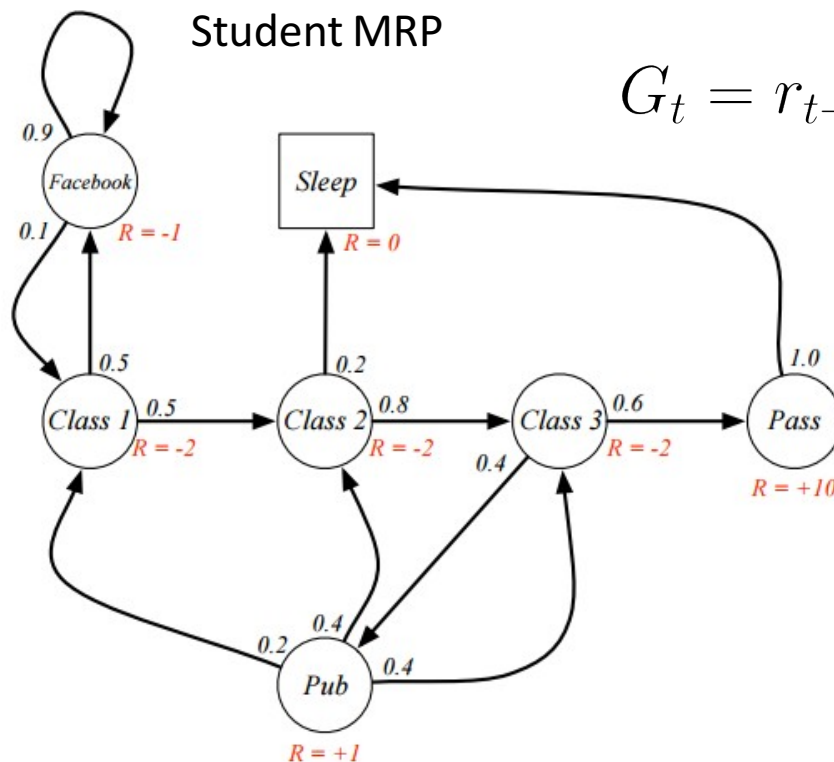
Sample episodes from $S_1=C1$

- C1 C2 C3 Pass Sleep
- C1 FB FB C1 C2 Sleep
- C1 C2 C3 Pub C2 C3 Pass Sleep
- C1 FB FB C1 C2 C3 Pub
- C1 FB FB FB C1 C2 C3 Pub C2 Sleep

Markov Reward Process (MRP)

Markov reward process is a Markov chain with values

- The return G_t is the total discounted reward from time-step t

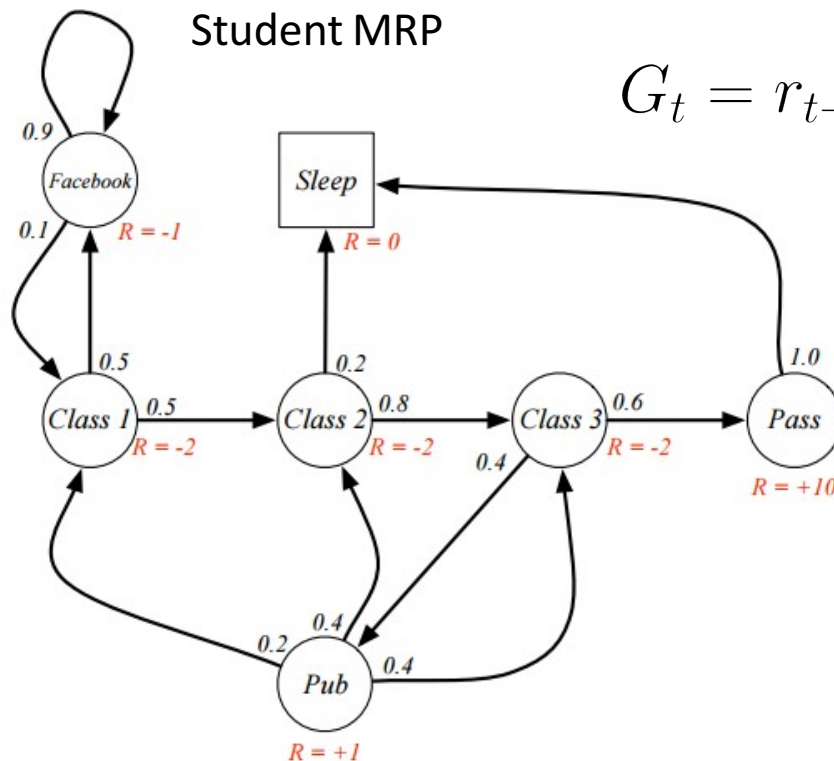


$$G_t = r_{t+1} + \gamma r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

Quick question: why do we need a discount factor in Markov Reward Process (MRP)?

Markov reward process is a Markov chain with values

- The return G_t is the total discounted reward from time-step t

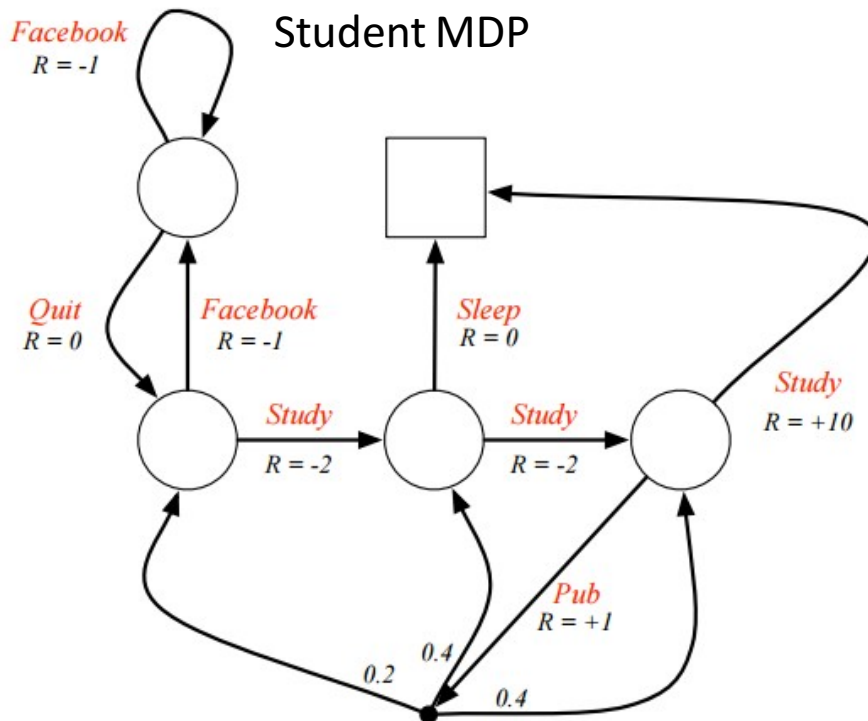


$$G_t = r_{t+1} + \gamma r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

Markov Decision Process (MDP)

Markov decision process is a MRP with decisions

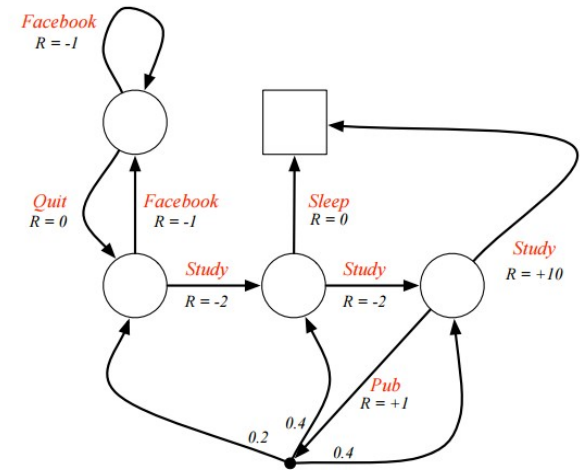
- It is an environment in which all states are Markov



Markov Decision Process (MDP)

- S : finite set of **states/observations** A : finite set of **actions**
- P : transition **probability** R : immediate **reward**
- γ : discount factor
- Goal is to choose **policy** π at time t that maximizes expected
- overall return:

$$\sum_{t'=t}^T \gamma^{t'-t} r_{t'}$$



Reinforcement Learning

Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

Reinforcement Learning

- Policy-Based
- Value-Based
- Model-Based

Problems within RL

- Learning and Planning
- Exploration and Exploitation

Major Components in an RL Agent

An RL agent may include one or more of these components

- **Policy**: agent's behavior function
- **Value function**: how good is each state and/or action
- **Model**: agent's representation of the environment

Policy

A policy is the agent's behavior

A policy maps from state to action

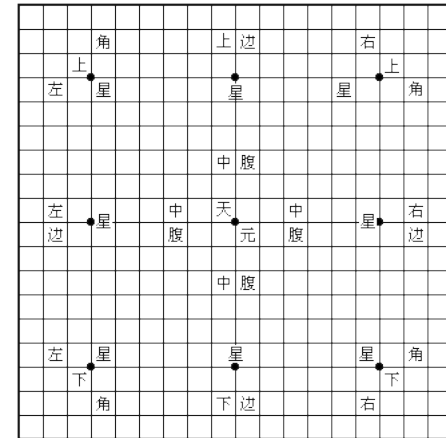
- Deterministic policy: $a = \pi(s)$
- Stochastic policy: $\pi(a) = P(a \mid s)$

Value Function

A value function is a prediction of future reward (with action a in state s)

Q.-value function gives expected total reward

- from state S and action A
- under policy π
- with discount factor γ



$$Q^{\pi}(s, a) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

Value functions decompose into a Bellman equation

$$Q^{\pi}(s, a) = \mathbb{E}_{s', a'}[r + \gamma Q^{\pi}(s', a') \mid s, a]$$

Optimal Value Function

An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$$

The optimal value function allows us act optimally

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

The optimal value informally maximizes over all decisions

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

Optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$$

Piazza Poll: Is Bellman equation a dynamic programming equation?

An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$$

The optimal value function allows us act optimally

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

The optimal value informally maximizes over all decisions

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

Optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$$

Model

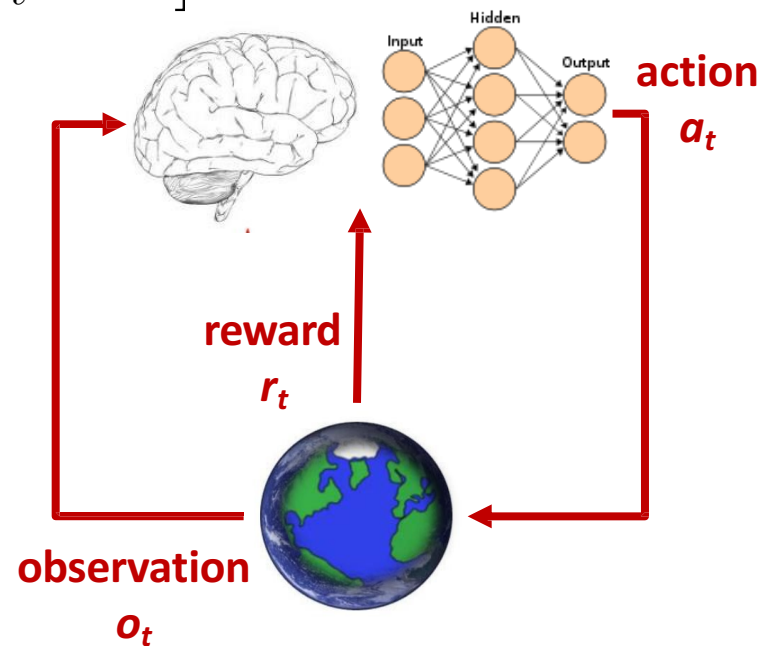
A model predicts what the environment will do next

- P predicts the next state

$$P_{ss'}^a = \mathbb{P}[s_{t+1} = s' \mid s_t = s, a_t = a]$$

- R predicts the next immediate reward

$$R_s^a = \mathbb{E}[r_{t+1} \mid s_t = s, a_t = a]$$



Reinforcement Learning Approach

Policy-based RL

- Search directly for optimal policy π^*

π^* is the policy achieving maximum future reward

Value-based RL

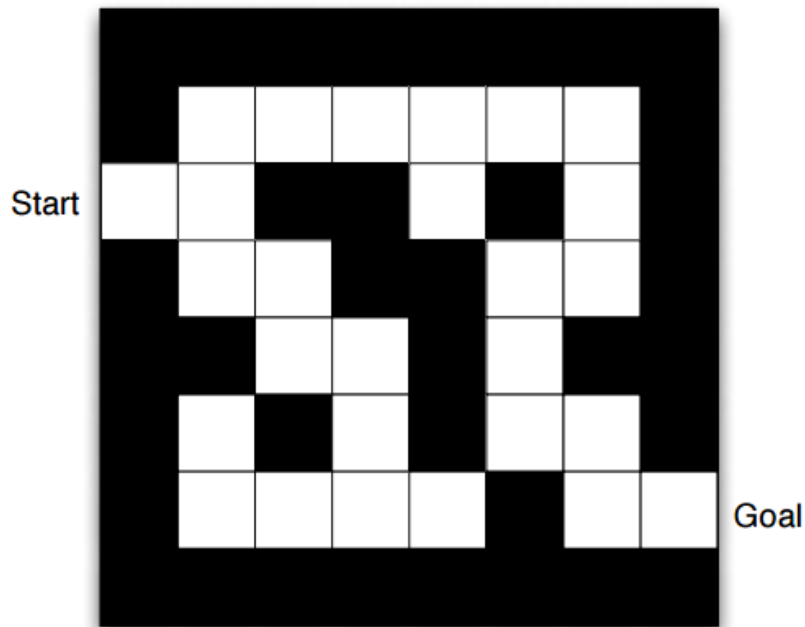
- Estimate the optimal value function $Q^*(s, a)$

$Q^*(s, a)$ is maximum value achievable under any policy

Model-based RL

- Build a model of the environment
- Plan (e.g. by lookahead) using model

Maze Example

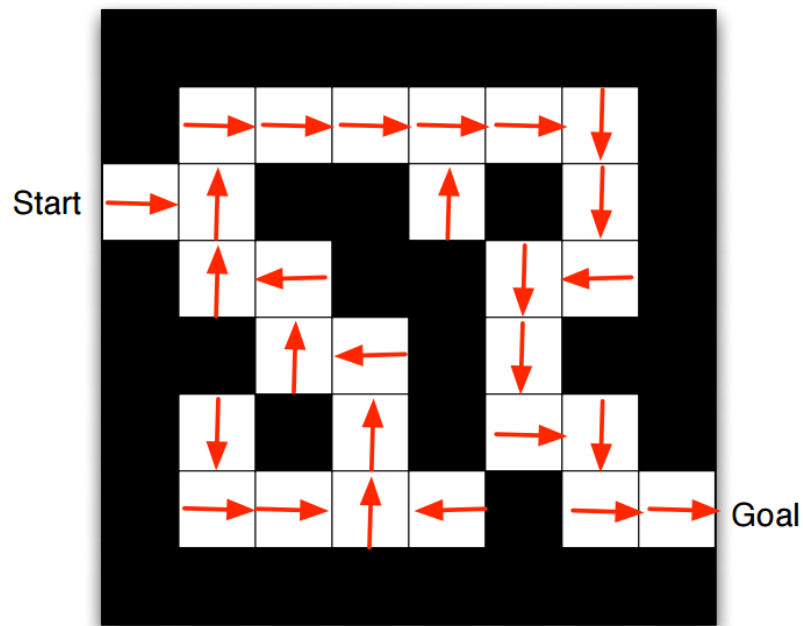


Rewards: -1 per time-step

Actions: N, E, S, W

States: agent's location

Maze Example: Policy



Rewards: -1 per time-step

Actions: N, E, S, W

States: agent's location

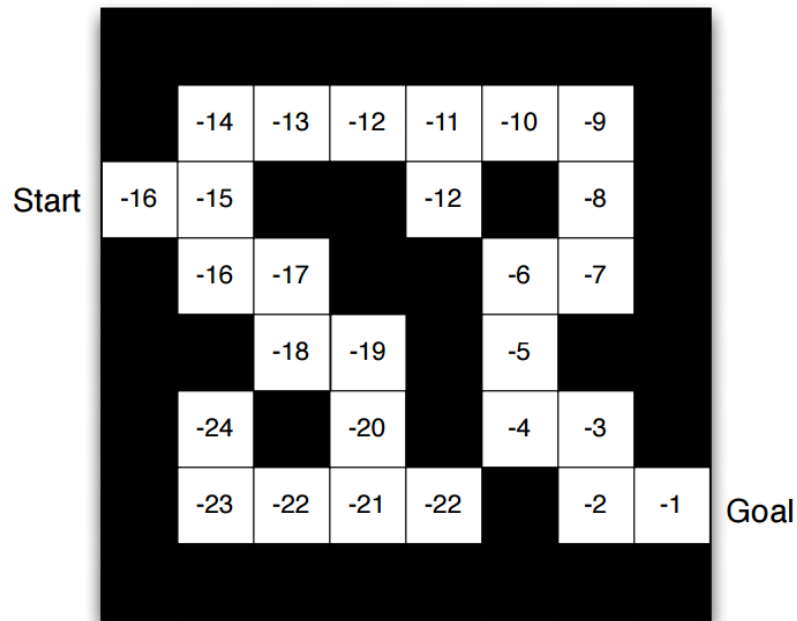
Arrows represent policy $\pi(s)$ for each state s

Maze Example: Value Function

Rewards: -1 per time-step

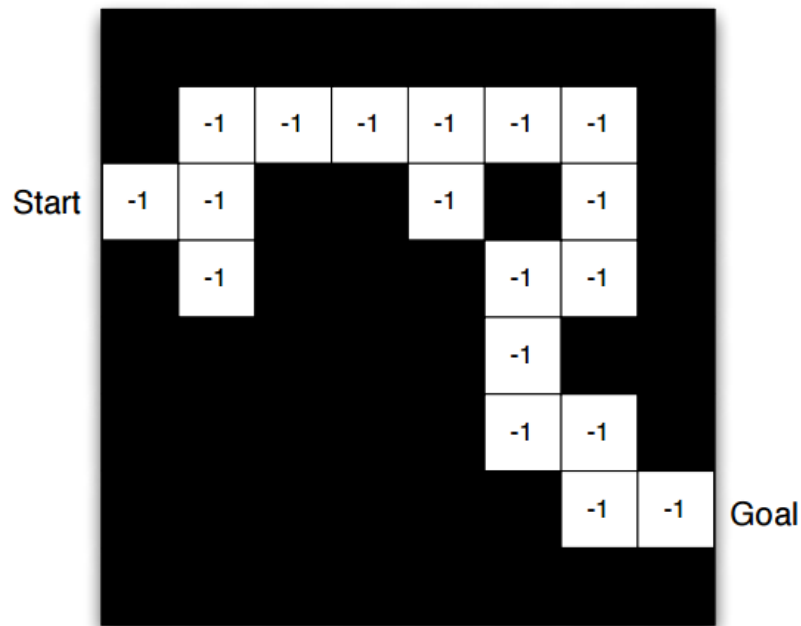
Actions: N, E, S, W

States: agent's location



Numbers represent value $Q_{\pi}(s)$ of each state s

Maze Example: Value Function



Rewards: -1 per time-step

Actions: N, E, S, W

States: agent's location

Grid layout represents transition model P

Numbers represent immediate reward R from each state s (same for all a)

Categorizing RL Agents

Model-Free

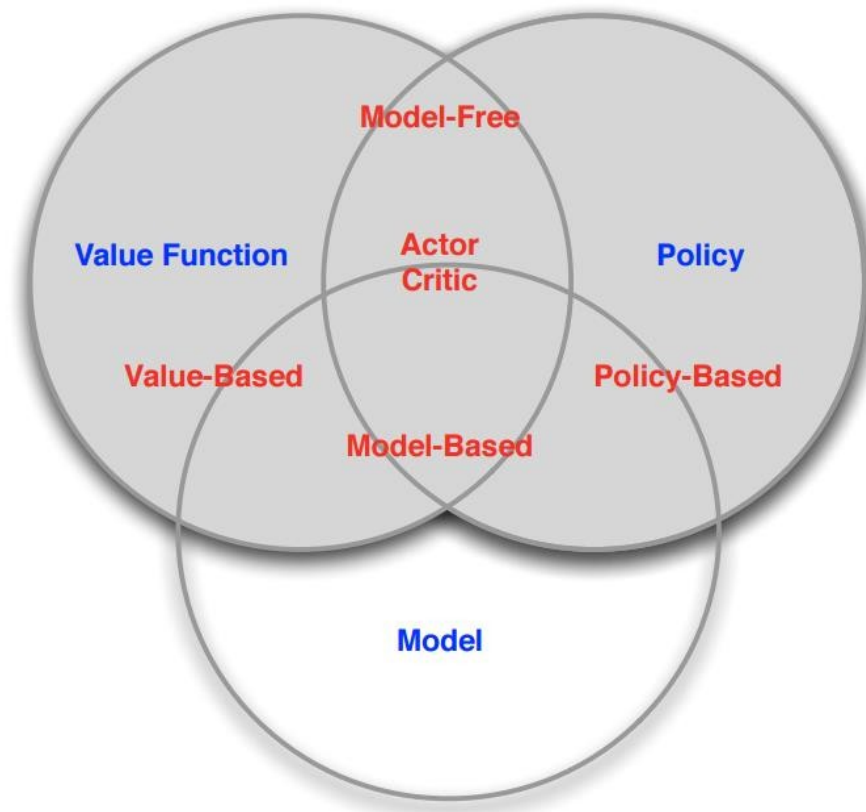
- Policy and/or Value Function
- No Model

Model-Based

- Policy and/or Value Function
- Model

- Value-Based
 - No Policy (implicit)
 - Value Function
- Policy-Based
 - Policy
 - No Value Function
- Actor-Critic
 - Policy
 - Value Function

RL Agent Taxonomy



Problems within RL

Outline

Machine Learning

- Supervised Learning v.s. Reinforcement Learning
- Reinforcement Learning v.s. Deep Learning

Introduction to Reinforcement Learning

- Agent and Environment
- Action, State, and Reward

Markov Decision Process

Reinforcement Learning

- Policy-Based
- Value-Based
- Model-Based

Problems within RL

- Learning and Planning
- Exploration and Exploitation

Learning and Planning

In sequential decision making

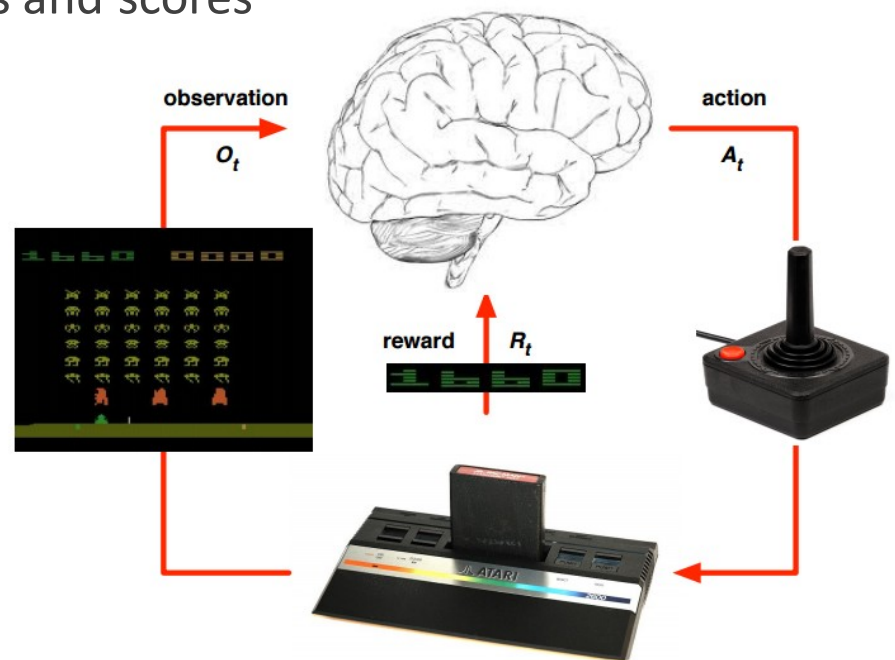
- Reinforcement learning
 - The environment is initially unknown
 - The agent interacts with the environment
 - The agent improves its policy
- Planning
 - A model of the environment is known
 - The agent performs computations with its model (w/o any external interaction)
 - The agent improves its policy (a.k.a. deliberation, reasoning, introspection, pondering, thought, search)

Atari Example: Reinforcement Learning

Rules of the game are unknown

Learn directly from interactive game-play

Pick actions on joystick, see pixels and scores



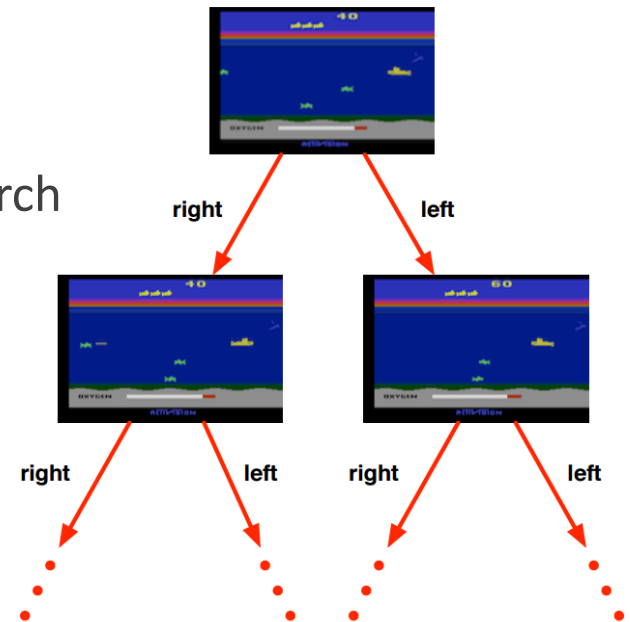
Atari Example: Planning

Rules of the game are known

Query emulator based on the perfect model inside agent's brain

- If I take action a from state s :
 - what would the next state be?
 - what would the score be?

Plan ahead to find optimal policy e.g. tree search



Exploration and Exploitation

Reinforcement learning is like **trial-and-error** learning

The agent should discover a good policy from the experience without losing too much reward along the way

When to try?

Exploration finds more information about the environment

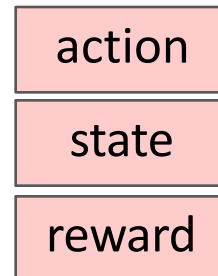
Exploitation exploits known information to maximize reward

It is usually important to explore as well as exploit

Concluding Remarks

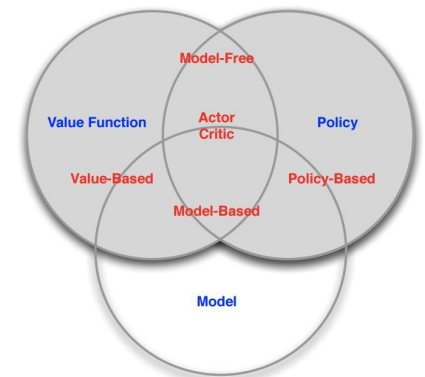
RL is a general purpose framework for **decision making** under interactions between *agent* and *environment*

- RL is for an *agent* with the capacity to *act*
- Each *action* influences the agent's future *state*
- Success is measured by a scalar *reward* signal
- Goal: *select actions to maximize future reward*



An RL agent may include one or more of these components

- **Policy**: agent's behavior function
- **Value function**: how good is each state and/or action
- **Model**: agent's representation of the environment



References

Course materials by David Silver: <http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html>

ICLR 2015 Tutorial: <http://www.iclr.cc/lib/exe/fetch.php?media=iclr2015:silver-iclr2015.pdf>

ICML 2016 Tutorial: http://icml.cc/2016/tutorials/deep_rl_tutorial.pdf