Lecture 1 : Intro

Reinforcement Learning aims to learn to
make good sequences of decisions

                                             of
Learn   to   make   good   sequence decisions
  |                      |        repeated interaction
  ↓                   learn to        with world
· don't know m         make
  advance           good sequences
who how the       of decision
world works

· fundamental challenge in artificial intelligence
and machine learning is learning to make
good decisions under uncertainty.

RL, behavior & Intelligence
_____

    Yael Niv
    Childhood : primitive brain & eye, swims around,
                attaches to a rock
    Adulthood : digests brain and sit7
    → Brain is helping guide decision [no more decisions, no need for brain?]

DeepMind Nature, 2015 (Atari Learning)

Robotics
Educational Games
☙ Health care

NLP, Vision

RL involves
_____

Optimization
Delayed consequences

Exploration
Generalization

Goal is to find an optimal way to make decisions
Or at least a very good strategy

Decisions now can impact things much later (saving for retirement)

Introduces two challenges (when planning and when learning)

Finding key of ~~Mucha~~ ~~leverage~~ Montezuma revenge

Exploration:
_____

• Learning about the world by making decisions

• Censored data (Reward is the only way)

• Decisions impact what we learn about

Policy is mapping from past experience to action
_____

Why we can't pre program it? (Not possible in big cases)

# Reinforcement Learning 2019 Stanford

## Lecture 1 Intro

$\boxed{RL}$

O - Optimization
G - Generalization
E - Exploration
D - Delayed Consequency

### AI Planning vs RL

| O | G | E | D |
|---|---|---|---|
| ✓ | ✓ | ✗ | ✓ |

### Superised ML vs RL

| O | G | E | D |
|---|---|---|---|
| ✓ | ✓ | ✗ | ✗ |

Learns from experience

### Unsupervise ML vs RL

| O | G | E | D |
|---|---|---|---|
| ✓ | ✓ | ✗ | ✗ |

### Imitation ~~Ree~~ Learning   [popularised by Andrew Ng]

| O | G | E | D |
|---|---|---|---|
| ✓ | ✓ | ✗ | ✓ |

Learns from experience of others.

Assumes input demos of good policies.

- Reduces RL to supervised learning

Benifits
- Great tools for supervised learning
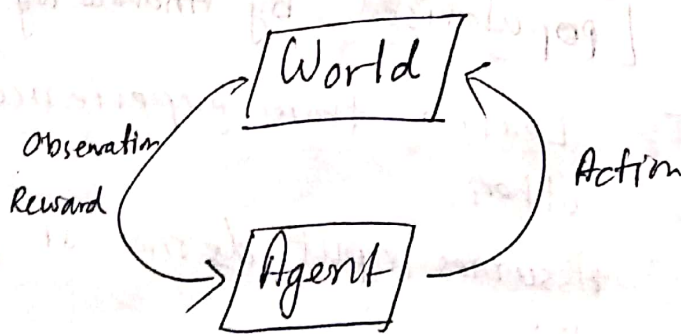- Avoids exploration problem

- Limitations
  - Can be expensive to capture
  - Limited by data collected.

Imitation ~~Real~~ Learning + RL promising

## How do we proceed?

- Explore the world
- Use experience to guide future decisions

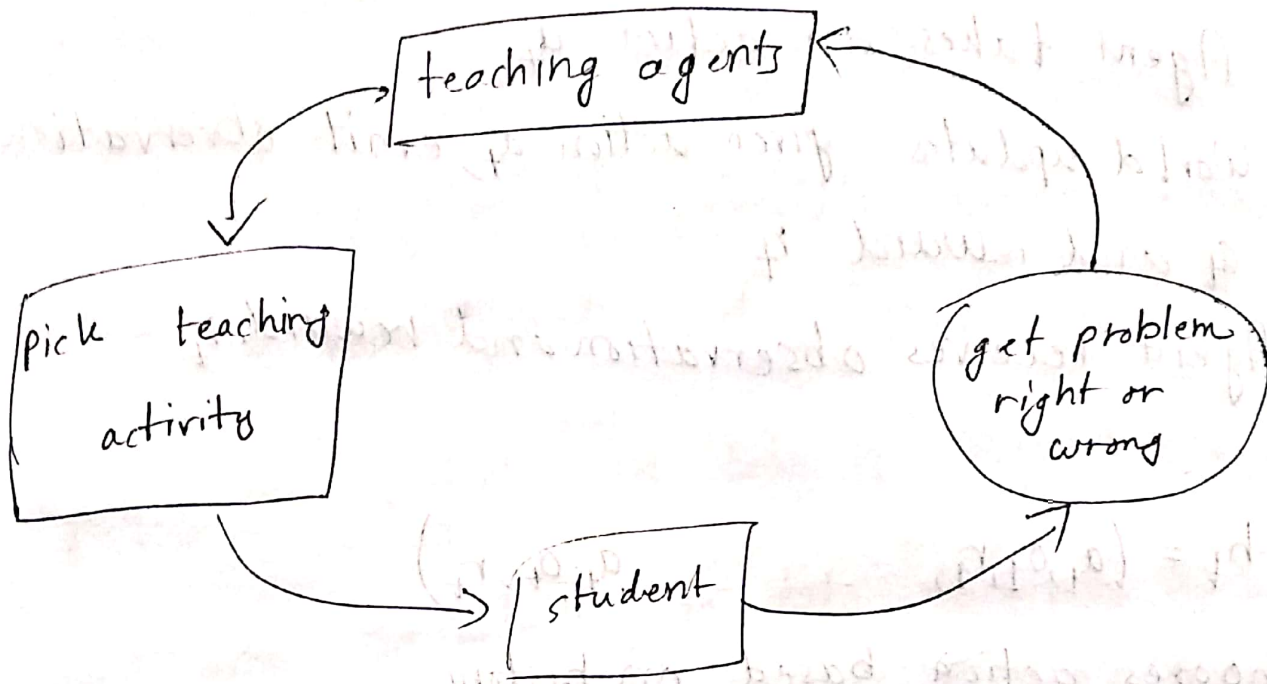- Introduction to sequential decision making under uncertainty



Example: Web Advertising (i)
Robot unloading (ii)
Dichwasher
Blood Pressure control (i) (ii)

Goal: select actions to maximize total expected future reward

(i) • May require balancing intermediate & long-term ~~resut~~ rewards

(ii) • May require strategic behavior to achieve high rewards

# Artificial Tutor:



+1 if the student get the problem right

-1 if they get it wrong

what activity would agent choose to get max $\Sigma$rewards?

Bev Wolf, 2000

The • Student initially doesn't know addition (easier) nor
substraction (harder)

    Reward hacking → give easy problems first
        then hard problems

    • Machine teaching

Each time step $t$:

- Agent takes an action $a_t$
- World updates given action $a_t$, emit observation $o_t$ and reward $r_t$
- Agent receives observation and reward. $r_t$

History $h_t = (a_1, o_1, r_1, \ldots a_t, o_t, r_t)$

Agent chooses action based on history

State is information & assumed to determine
- what happens next. $S_t = (h_t)$

**World State**
- This is the true state of the world used to determine how world generates next observation and reward

- often hidden and unknown to agent
- Even if & known may contain information not needed by agent.

**Agent State** What the agent/algorithm uses to make decisions about how to act $s_t = f(h_t)$
Generally including meta info. like state of algo. or decision process

# Markov Assumption

- Information stat: sufficient static of history
- State $s_t$ is Markov if and only if

$$P(s_{t+1}|s_t, a_t) = P(s_{t+1}|h_t, a_t)$$

- Future is independent of past given present

## Hypertension control:

let state be current blood pressure, and action be whether to take decisions & medications or not. Is the system markov? (No)

## Website shopping:

state is current product viewed by customer and action is what other product to recomment. Is the system Markov? (No)

## Why is it popular?

- Can always be satisfied
  [setting state as history always Markov; $s_t = h_t$]

- In practive $s_t = a_t$

State representation has a big implication for
- Computational complexity
- Data required
- Resulting performance

# Full Observability (MDP)

$$S_t = O_t$$

## Partial Observability (POMDP)

Agent state is not the same as world state

Agent constructs its own state

· Use history $S_t = h_t$ or beliefs of world state

Poker healthcare

# Types of Sequential Decision Processes

Bandits : actions have no influence on next observation

No delayed rewards.

## How world changes

Deterministic (single observation and reward)

Stochastic (many " " " )

## RL algo components

Model → representation of how the world changes

Policy → function mapping agent's states to action
$$\pi(s) = a \quad \pi(a|s) = Pr(o_{t=a}|s_{t=a})$$

Value function : Future rewards from being in a state and/or action when following a particular policy

Value function : Expected discounted sum of
future rewards, under ~~the~~ a particular ~~poli~~ policy π

$$v^{\pi}(s_t^d = s) = E_{\pi}\left[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} \cdots - - \quad s_t = s\right]$$

Discount factor $\gamma$ weights immediate vs future
rewards
$$0 \longleftrightarrow 1$$

~~Can be quantifies goodne~~

Can be used to quantify goodness/badness
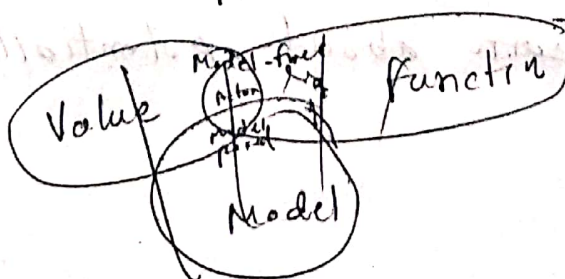Decide how to act ~~compatly~~ comparing policies.

Types of RL agents
_____

• Model-based • Explicit: Mode
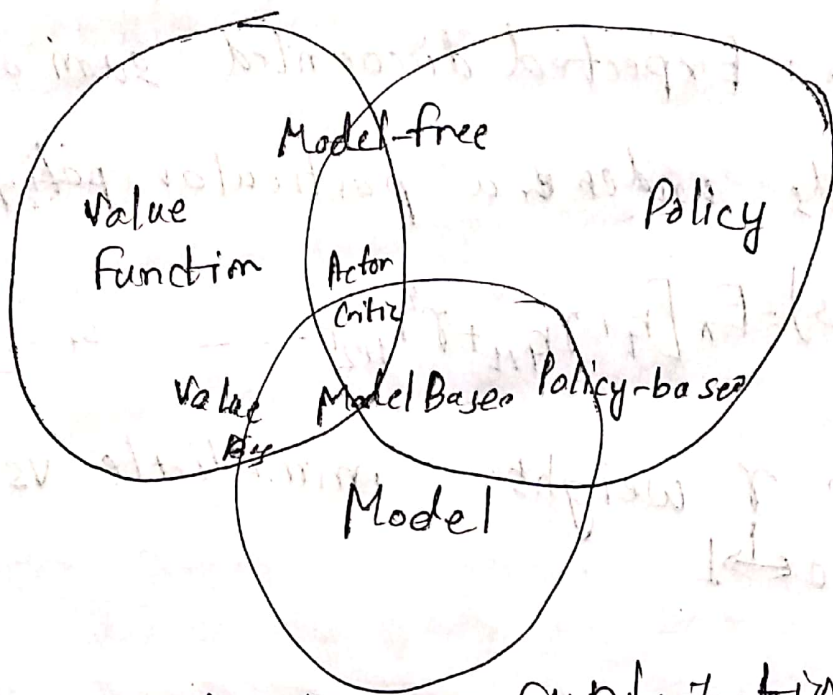                • May or maynot have policy and/or
                  value funcs
• Model-free    • Explicit : Value function    and/or ~~a~~ policy
                                                        function

                • No model


Value  Model-free  function
        Model

Value Function

Model-free

Policy

Actor Critic

Value By

Model Based Policy-based

Model

Exploration an ~~Exploratin~~ Exploitation

Exploration and Exploitation:

→ choosing actions that are expected to yield good reward given past experience

→ trying new things that might enable the agent to make better decisions in the future

• May have to sacrifice reward in order to explore & learn about potentially better policy