

# CS 440/ECE448: Model-Based Reinforcement Learning

Mark Hasegawa-Johnson

These slides are in the public domain.



By Nicolas P. Rougier - Own work, CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=29327040>

# Outline

- Reinforcement learning
- Model-based learning
- Exploration vs. Exploitation

# Review: Markov Decision Process

- MDP defined by states, actions, transition model, reward function
- The “solution” to an MDP is the policy: what do you do when you’re in any given state
- The Bellman equation tells the utility of any given state, and incidentally, also tells you the optimum policy. The Bellman equation is  $N$  nonlinear equations in  $N$  unknowns (the policy), therefore it can’t be solved in closed form.
- Value iteration:
  - At the beginning of the  $(i + 1)^{\text{st}}$  iteration, each state’s value is based on looking ahead  $i$  steps in time
  - ... so finding the best action = optimize based on  $(i + 1)$ -step lookahead
- Policy iteration:
  - Find the utilities that result from the current policy,
  - Improve the current policy

# Reinforcement learning: Basic scheme

But what if you don't know  $P(s'|s, a)$  or  $r(s)$ ?

Answer: “learning by doing” (a.k.a. reinforcement learning).

In each time step:

- Take some action
- Observe the outcome of the action: successor state and reward
- Update some internal representation of the environment and policy

# Key problems

## 1. What should you learn?

- Model-based learning:  $P(s'|s, a)$  and  $r(s)$  are estimated using a neural network or probability table, then use value iteration or policy iteration to find the best policy
- Model-free learning:
  - Policy learning:  $\pi(s)$ , the policy, is directly estimated using a neural network or a table
  - Q-learning:  $Q(s, a)$ , the quality of action  $a$ , is estimated using a neural network or a table of numbers, and directly specifies the best action

## Key problems, 2. In which order should you study the states?

- Real-time learning
  - In state  $s_t$ , try action  $a_t$ , see what reward  $r_t$  state  $s_{t+1}$  results, and immediately update your estimates of  $r(s_t)$  and  $P(s_{t+1}|s_t, a_t)$
- Experience replay buffer
  - In state  $s_t$ , try action  $a_t$ , see what reward  $r_t$  state  $s_{t+1}$  results, and store the tuple  $(s_t, a_t, r_t, s_{t+1})$  in an experience replay buffer
  - When the experience replay buffer is full, learn by drawing samples from it according to some criterion that optimizes the rate at which you learn

# Key problems, 3. Which actions should you perform while learning?

- Real-time Learning: Exploration vs. Exploitation
  - Exploration: try actions at random, to see what happens
  - Exploitation: try to act optimally (to maximize value)
- Experience replay buffer: On-policy vs. Off-policy learning:
  - On-policy: Learn using tuples  $(s_t, a_t, r_t, s_{t+1})$  where  $a_t = \pi(s_t)$  matches your current beliefs about the best policy
  - Off-policy: Learn using tuples  $(s_t, a_t, r_t, s_{t+1})$  where the action was chosen according to some old policy that you no longer believe is the best, but the tuple is in your experience replay buffer, so you learn from it anyway

# Outline

- Reinforcement learning
- Model-based learning
- Exploration vs. Exploitation

# Model-based real-time learning

- Take an action
- Observe the state and the reward that result
- Update your current estimates of  $P(s'|s, a)$  and  $r(s)$

# Example of model-based reinforcement learning: Theseus the Mouse



[Claude Shannon and Theseus the Mouse](#). Public domain image, Bell Labs.

[https://www.youtube.com/watch?v=9\\_AEVQ\\_p74](https://www.youtube.com/watch?v=9_AEVQ_p74)

# Model-based reinforcement learning: Theseus' strategy

- At each position in the maze ( $s$ ),
  - For every possible action  $a \in \{\text{Forward, Left, Right, Back}\}$ :
    - If the action succeeded in changing the state ( $s' \neq s$ ), then set  $P(s'|s, a) = 1$
    - If not, set  $P(s'|s, a) = 0$  for all  $s' \neq s$

Once you've learned the maze, then compute the best policy ( $\pi(s)$ ) using value iteration or policy iteration.

Learning phase:

- At each position in the maze ( $s$ ),
  - For every possible action  $a \in \{\text{Forward, Left, Right, Back}\}$ :
    - If the action succeeded in changing the state ( $s' \neq s$ ), then set  $P(s'|s, a) = 1$
    - If not, set  $P(s'|s, a) = 0$  for all  $s' \neq s$

Once you've learned the maze, then compute the best policy ( $\pi(s)$ ) using Value Iteration.

- If  $P(s'|s, a) \in (0,1)$ , Value Iteration = BFS

# Learning non-trivial probabilities

- Let's keep a table of numbers,  $N(s, a, s')$ , telling how many times action  $a$  in state  $s$  led to next-state  $s'$
- At time  $t$ , in state  $s_t$ :
  - choose action  $a_t$
  - observe  $r_t$  and  $s_{t+1}$
  - learn:

$$N(s_t, a_t, s_{t+1}) += 1$$
$$P(s_{t+1}|s_t, a_t) = \frac{N(s_t, a_t, s_{t+1}) + k}{\sum_{s' \in \mathcal{S}} N(s_t, a_t, s') + k|\mathcal{S}|}$$

# Outline

- Reinforcement learning
- Model-based learning
- **Exploration vs. Exploitation**

# Exploration vs. Exploitation

- **Exploration:** take a new action with unknown consequences
  - Pros:
    - Get a more accurate model of the environment
    - Discover higher-reward states than the ones found so far
  - Cons:
    - When you're exploring, you're not maximizing your utility
    - Something bad might happen
- **Exploitation:** go with the best strategy found so far
  - Pros:
    - Maximize reward as reflected in the current utility estimates
    - Avoid bad stuff
  - Cons:
    - Might also prevent you from discovering the true optimal strategy

“Search represents a core feature of cognition:”  
[Exploration versus exploitation in space, mind, and society.](#)

# How to trade off exploration vs. exploitation

**Epsilon-first strategy** ( $\epsilon$  is an integer): when you reach state  $s$ , check how many times you've tested each of its available actions.

- **Explore for the first  $\epsilon$  trials**: If the least-explored action has been tested fewer than  $\epsilon$  times, then perform that action.
- **Exploit thereafter**: Once you've finished exploring, start exploiting (work to maximize your personal utility).
- Example: Theseus.

**Epsilon-greedy strategy** ( $0 < \epsilon < 1$ ): in every state, every time, forever,

- **With probability  $\epsilon$ , Explore**: choose any action, uniformly at random.
- **With probability  $(1 - \epsilon)$ , Exploit**: choose the action with the highest expected utility, according to your current estimates.
- Guarantee:  $P(s'|s, a)$  converges to its true value as #trials  $\rightarrow \infty$ .

# The epsilon-first strategy

As you wander through the maze, you reach some state,  $s$ .

- If there is any action,  $a$ , for which  $N(s, a) < \epsilon$ , then **explore** (= try the action, to see what it does).
- If not, then **exploit** your knowledge (choose the action that, according to your model, will lead to the highest utility).



Claude Shannon and Theseus the Mouse. Public domain image, Bell Labs.

# The epsilon-first strategy

- Remember, we're learning transition probabilities as:

$$P(s_{t+1}|s_t, a_t) \approx \frac{N(s_t, a_t, s_{t+1})}{N(s_t, a_t)}$$

- If we try every action  $N(s_t, a_t) = \epsilon$  times, then our probability estimates are correct to a resolution of about  $1/\epsilon$ :

$$P(s_{t+1}|s_t, a_t) \approx \frac{N(s_t, a_t, s_{t+1})}{\epsilon}$$

- For example, if you want your estimates to be correct to the nearest 10%, then you should try every action at least 10 times.



Claude Shannon and Theseus the Mouse. Public domain image, Bell Labs.

# The epsilon-greedy strategy

Regardless of how few times or how many times you've been in state  $s$ : generate a uniform random number,  $z \in (0,1)$ .

- If  $z \leq \epsilon$ , then **explore**. Choose an action,  $a$ , uniformly at random, and try it. See what  $s'$  results. Increment  $N(s, a)$  and  $N(s, a, s')$ .
  - This happens with probability  $\epsilon$ .
- If  $z > \epsilon$ , then **exploit**. Use value iteration or policy iteration to figure out the best action in the current state, then do that action.
  - This happens with probability  $1 - \epsilon$ .

# Quiz

Try the quiz!

# Compare: Epsilon-first and Epsilon-greedy

$$\text{For both: } P(s'|s, a) \approx \frac{N(s, a, s')}{N(s, a)}$$

## Advantages of Epsilon-first:

- In the beginning, when  $P(s'|s, a)$  is still inaccurate, we just try things at random (explore).
- We can choose the level of precision that's "enough" for us. When  $P(s'|s, a)$  reaches that point, we stop exploring, and instead, we focus on getting the biggest rewards possible (exploit).

## Advantages of Epsilon-greedy:

- Gradually, over a series of many experiments,  $N(s, a) \rightarrow \infty$
- Therefore, as the number of experiments gets large, our estimates of the probabilities converge to their true values

# Summary

- Reinforcement learning
- Model-based learning

$$P(s_{t+1}|s_t, a_t) = \frac{N(s_t, a_t, s_{t+1}) + k}{\sum_{s' \in \mathcal{S}} N(s_t, a_t, s') + k|\mathcal{S}|}$$

- Exploration vs. Exploitation
  - Epsilon-first: explore every action at least  $\epsilon$  times
  - Epsilon-greedy: explore at random with probability  $\epsilon$