

# Artificial Intelligence Foundation – JC3001

Lecture 44: Reinforcement Learning - I

**Prof. Aladdin Ayesh** (aladdin.ayesh@abdn.ac.uk)

**Dr. Binod Bhattarai** (binod.bhattarai@abdn.ac.uk)

**Dr. Gideon Ogunniye**, (g.ogunniye@abdn.ac.uk)

October 2025

Material adapted from:  
Russell and Norvig (AIMA Book): Chapter 22  
Sutton and Barto (Reinforcement Learning: An Introduction 2nd ed.)  
David Silver (UCL)  
Michael Littman (Brown University) and Charles Isbell (GA Tech)

## Course Progression

- Part 1: Introduction
  - ① Introduction to AI ✓
  - ② Agents ✓
- Part 2: Problem-solving
  - ① Search 1: Uninformed Search ✓
  - ② Search 2: Heuristic Search ✓
  - ③ Search 3: Local Search ✓
  - ④ Search 4: Adversarial Search ✓
- Part 3: Reasoning and Uncertainty
  - ① Reasoning 1: Constraint Satisfaction ✓
  - ② Reasoning 2: Logic and Inference ✓
  - ③ Probabilistic Reasoning 1: BNs ✓
  - ④ Probabilistic Reasoning 2: HMMs ✓
- Part 4: Planning
  - ① Planning 1: Intro and Formalism ✓
  - ② Planning 2: Algorithms & Heuristics ✓
  - ③ Planning 3: Hierarchical Planning ✓
  - ④ Planning 4: Stochastic Planning ✓
- Part 5: Learning
  - ① Learning 1: Intro to ML ✓
  - ② Learning 2: Regression ✓
  - ③ Learning 3: Neural Networks ✓
  - ④ **Learning 4: Reinforcement Learning**
- Part 6: Conclusion
  - ① Ethical Issues in AI
  - ② Conclusions and Discussion

# Objectives

- Bandit Problems
- Reinforcement Learning based Agents
- Tabular Reinforcement Learning
- Function Generalization



# Outline

## 1 Recap and Motivation

► Recap and Motivation

► Multi-armed Bandits

► Reinforcement Learning based Agents

# Motivation

## 1 Recap and Motivation

	1	2	3	4
a				+1
b				-1
c				

- Recall the grid world example we saw for MDPs

# Motivation

## 1 Recap and Motivation

	1	2	3	4
a				+1
b				-1
c				

- Recall the grid world example we saw for MDPs
- We learned how to calculate the optimal policy using transition **probabilities** and **rewards**

# Motivation

## 1 Recap and Motivation

	1	2	3	4
a				+1
b				-1
c				

- Recall the grid world example we saw for MDPs
- We learned how to calculate the optimal policy using transition **probabilities** and **rewards**
- But what if we do not know the rewards, probabilities or both?



# Motivation

## 1 Recap and Motivation

	1	2	3	4
a				+1
b				-1
c				

- Recall the grid world example we saw for MDPs
- We learned how to calculate the optimal policy using transition **probabilities** and **rewards**
- But what if we do not know the rewards, probabilities or both?
- We should be able to explore the environment and discover the best policy

# Learning Recap

## 1 Recap and Motivation

Main methods of learning:

- Supervised learning:

From  $(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots (\vec{x}^{(m)}, y^{(m)})$ , derive a function  $f(X^{(m)}) = Y^{(m)}$

# Learning Recap

## 1 Recap and Motivation

Main methods of learning:

- Supervised learning:

From  $(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots, (\vec{x}^{(m)}, y^{(m)})$ , derive a function  $f(X^{(m)}) = Y^{(m)}$

- Unsupervised learning:

From  $\vec{x}^{(1)}, \vec{x}^{(2)}, \vec{x}^{(m)}, \dots, \vec{x}^{(m)}$ , identify classes  $x_c$  and derive a function to calculate  $P(X = x_c)$

# Learning Recap

## 1 Recap and Motivation

Main methods of learning:

- Supervised learning:

From  $(\vec{x}^{(1)}, y^{(1)}), (\vec{x}^{(2)}, y^{(2)}), \dots, (\vec{x}^{(m)}, y^{(m)})$ , derive a function  $f(X^{(m)}) = Y^{(m)}$

- Unsupervised learning:

From  $\vec{x}^{(1)}, \vec{x}^{(2)}, \vec{x}^{(m)}, \dots, \vec{x}^{(m)}$ , identify classes  $x_c$  and derive a function to calculate  $P(X = x_c)$

- Reinforcement learning:

From a sequence of state transitions (i.e. states and actions) and occasional rewards  $[s \rightarrow^r a, s \rightarrow^r a, \dots]$ , discover the optimal policy  $\pi^*(s)$  for this environment.

# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.    Unsup.    Reinf.

Spam Filtering

Speech Recognition

Starlight classification

Playing SF $\Xi$  against an unknown player

Navigating an unknown environment

Elevator controller

# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.    Unsup.    Reinf.

X

Spam Filtering

Speech Recognition

Starlight classification

Playing SF $\Xi$  against an unknown player

Navigating an unknown environment

Elevator controller

# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.    Unsup.    Reinf.

X

X

Spam Filtering

Speech Recognition

Starlight classification

Playing SF $\Xi$  against an unknown player

Navigating an unknown environment

Elevator controller

# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.    Unsup.    Reinf.

X

X

X

Spam Filtering

Speech Recognition

Starlight classification

Playing SF $\Xi$  against an unknown player

Navigating an unknown environment

Elevator controller



# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.	Unsup.	Reinf.	
X			Spam Filtering
X			Speech Recognition
	X		Starlight classification
		X	Playing SF $\Xi$ against an unknown player
			Navigating an unknown environment
			Elevator controller

# Machine Learning Problems

## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.	Unsup.	Reinf.	
X			Spam Filtering
X			Speech Recognition
	X		Starlight classification
		X	Playing SF $\Xi$ against an unknown player
		X	Navigating an unknown environment
			Elevator controller

# Machine Learning Problems

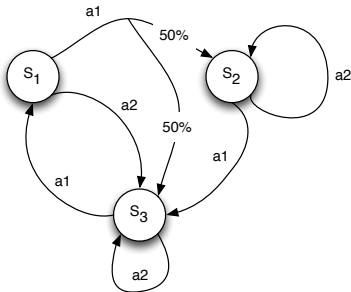
## 1 Recap and Motivation

Which of the previous methods of learning would be more suitable to solve the following problems:

Sup.	Unsup.	Reinf.	
X			Spam Filtering
X			Speech Recognition
	X		Starlight classification
		X	Playing SF $\Xi$ against an unknown player
		X	Navigating an unknown environment
		X	Elevator controller

# Markov Decision Process (MDP)

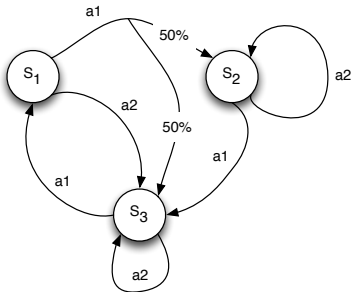
## 1 Recap and Motivation



- An MDP is defined in terms of
  - 1 An initial state  $S_0$
  - 2 A transition model  $T(s, a, s') = P(s'|a, s)$  (Markovian)
  - 3 A reward function  $R(s, a, s')$
- A solution to a MDP must specify what the agent should do for any state. Such a solution is called a policy

# Markov Decision Process (MDP)

## 1 Recap and Motivation



- An MDP is defined in terms of
  - 1 An initial state  $S_0$
  - 2 A transition model  $T(s, a, s') = P(s'|a, s)$  (Markovian)
  - 3 A reward function  $R(s, a, s')$  — sometimes expressed as  $R(s')$
- A solution to a MDP must specify what the agent should do for any state. Such a solution is called a policy

# Finding the Optimal Policy

## 1 Recap and Motivation

$$V(s) \leftarrow \left[ \max_a \gamma \sum_{s'} P(s'|s, a) * V(s') \right] + R(s)$$

# Finding the Optimal Policy

## 1 Recap and Motivation

$$V(s) \leftarrow \left[ \max_a \gamma \sum_{s'} P(s'|s, a) * V(s') \right] + R(s)$$

$$\pi^*(s) = \arg \max_a \sum_{s'} P(s'|s, a) * V(s')$$



# Outline

## 2 Multi-armed Bandits

► Recap and Motivation

► **Multi-armed Bandits**

► Reinforcement Learning based Agents



# k-Armed Bandit Problem

## 2 Multi-armed Bandits

### k-armed Bandit Problem

- Repeated choices among  $k$  different options or actions
- Each choice has a reward drawn from a stationary probability distribution
- Goal is to maximize rewards over a period of time



Thus, the value is the expected reward:

$$\begin{aligned} q_*(a) &\doteq \mathbb{E}[R_t \mid A_t = a] \forall a \in \{1, \dots, k\} \\ &= \sum_r p(r \mid a)r \end{aligned}$$

And to maximize it over time

$$\arg \max_a q_*(a)$$

# Computing the expected reward

## 2 Multi-armed Bandits

Consider a doctor with 3 possible treatments to offer a patient.

- Each treatment can either cure the patient ( $r = 1$ ), do nothing ( $r = 0$ ), or make the patient sicker ( $r = -1$ ), with probabilities ( $p(r \mid drug)$ ) as follows:

$$d_1 \quad \text{📺} \quad p(r = 1 \mid d = d_1) = 0.1, p(r = 0 \mid d = d_1) = 0.8, \text{ and } p(r = -1 \mid d = d_1) = 0.1$$

$$d_2 \quad \text{📺} \quad p(r = 1 \mid d = d_2) = 0.3, p(r = 0 \mid d = d_2) = 0.1, \text{ and } p(r = -1 \mid d = d_2) = 0.6$$

$$d_3 \quad \text{📺} \quad p(r = 1 \mid d = d_3) = 0.3, p(r = 0 \mid d = d_3) = 0.5, \text{ and } p(r = -1 \mid d = d_3) = 0.2$$

# Computing the expected reward

## 2 Multi-armed Bandits

Consider a doctor with 3 possible treatments to offer a patient.

- Each treatment can either cure the patient ( $r = 1$ ), do nothing ( $r = 0$ ), or make the patient sicker ( $r = -1$ ), with probabilities ( $p(r \mid drug)$ ) as follows:

$$d_1 \quad \text{📺} \quad p(r = 1 \mid d = d_1) = 0.1, p(r = 0 \mid d = d_1) = 0.8, \text{ and } p(r = -1 \mid d = d_1) = 0.1$$

$$d_2 \quad \text{📺} \quad p(r = 1 \mid d = d_2) = 0.3, p(r = 0 \mid d = d_2) = 0.1, \text{ and } p(r = -1 \mid d = d_2) = 0.6$$

$$d_3 \quad \text{📺} \quad p(r = 1 \mid d = d_3) = 0.3, p(r = 0 \mid d = d_3) = 0.5, \text{ and } p(r = -1 \mid d = d_3) = 0.2$$

Compute  $q_*$  for all drug options and help the doctor choose:

$$d_1 \quad \text{📺} \quad q(d = d_1) = (0.1 * 1) + (0.8 * 0) + (0.1 * -1) = 0$$

$$d_2 \quad \text{📺} \quad q(d = d_2) = (0.3 * 1) + (0.1 * 0) + (0.6 * -1) = -0.3$$

$$d_3 \quad \text{📺} \quad q(d = d_4) = (0.3 * 1) + (0.5 * 0) + (0.2 * -1) = 0.1$$

# Learning Action Values

## 2 Multi-armed Bandits

### Sample-Average Method

$$\begin{aligned} Q_t(a) &\doteq \frac{\text{sum of rewards for } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} \\ &= \frac{\sum_{i=1}^{t-1} R_i}{t-1} \end{aligned}$$

## Example

### 2 Multi-armed Bandits

$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$



$$Q_0(\text{slot machine}) = 0$$



$$Q_0(\text{slot machine}) = 0$$



$$Q_0(\text{slot machine}) = 0$$

## Example

### 2 Multi-armed Bandits

$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$



$$1$$

$$Q_1(\text{slot machine}) = 1$$



$$0$$

$$Q_1(\text{slot machine with red}) = 0$$



$$0$$

$$Q_1(\text{slot machine with green}) = 0$$



## Example

### 2 Multi-armed Bandits

$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$



1

0

$$Q_2(\text{slot machine}) = 0.5$$



0

0

$$Q_2(\text{slot machine}) = 0$$



0

0

$$Q_2(\text{slot machine}) = 0$$



## Example

### 2 Multi-armed Bandits

$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$



1

0

0

$$Q_3(\text{pill}) = 0.5$$



0

0

-1

$$Q_3(\text{pill}) = -1$$



0

0

0

$$Q_3(\text{pill}) = 0$$





## Example

### 2 Multi-armed Bandits

$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$



1

0

0

0

$$Q_4(\text{pill}) = 0.5$$



0

0

-1

0

$$Q_4(\text{red pill}) = -1$$



0

0

0

1

$$Q_4(\text{green pill}) = 1$$



# Example

## 2 Multi-armed Bandits



$$\frac{\sum_{i=1}^{t-1} R_i}{t-1}$$

1

0

0

0

0

$$Q_5(\text{white pill}) = 0.5$$

0

0

-1

0

0

$$Q_5(\text{red pill}) = -0.5$$

-0.5

0

0

0

1

0

$$= Q_5(\text{green pill}) = 1$$



# Incremental Update Rule

## 2 Multi-armed Bandits

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i$$

# Incremental Update Rule

## 2 Multi-armed Bandits

$$\begin{aligned} Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\ &= \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i) \end{aligned}$$

# Incremental Update Rule

## 2 Multi-armed Bandits

$$\begin{aligned}Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\&= \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i) \\&= \frac{1}{n} (R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i)\end{aligned}$$

# Incremental Update Rule

## 2 Multi-armed Bandits

$$\begin{aligned} Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\ &= \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i) \\ &= \frac{1}{n} (R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i) \end{aligned}$$
$$= \frac{1}{n} (R_n + (n-1)Q_n)$$

Recall:  $Q_n = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i$

# Incremental Update Rule

## 2 Multi-armed Bandits

$$\begin{aligned}Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\&= \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\&= \frac{1}{n} \left( R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\&= \frac{1}{n} (R_n + (n-1)Q_n) \\&= \frac{1}{n} (R_n + nQ_n - Q_n)\end{aligned}$$

$$\text{Recall: } Q_n = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i$$

# Incremental Update Rule

## 2 Multi-armed Bandits

$$\begin{aligned}Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\&= \frac{1}{n} (R_n + \sum_{i=1}^{n-1} R_i) \\&= \frac{1}{n} (R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i)\end{aligned}$$

$$\begin{aligned}&= \frac{1}{n} (R_n + (n-1)Q_n) \\&= \frac{1}{n} (R_n + nQ_n - Q_n) \\&= Q_n + \frac{1}{n} (R_n - Q_n)\end{aligned}$$

$$\text{Recall: } Q_n = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i$$



# Incremental Update Rule

## 2 Multi-armed Bandits

We will see this form of update throughout this lecture:

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

$$Q_{n+1} = Q_n + \frac{1}{n}(R_n - Q_n)$$

$$Q_{n+1} = Q_n + \alpha_n(R_n - Q_n)$$



# Outline

## 3 Reinforcement Learning based Agents

- ▶ Recap and Motivation
- ▶ Multi-armed Bandits
- ▶ Reinforcement Learning based Agents

# Reinforcement Learning based Agents

## 3 Reinforcement Learning based Agents

Agents execute **trials** in the environment, using some policy  $\pi$ .

Typical trials look like:

$$(1, 1)_{-.04} \rightsquigarrow (1, 2)_{-.04} \rightsquigarrow (1, 3)_{-.04} \rightsquigarrow (1, 2)_{-.04} \rightsquigarrow (1, 3)_{-.04} \rightsquigarrow (2, 3)_{-.04} \rightsquigarrow (3, 3)_{-.04} \rightsquigarrow (4, 3)_{+1}$$

$$(1, 1)_{-.04} \rightsquigarrow (1, 2)_{-.04} \rightsquigarrow (1, 3)_{-.04} \rightsquigarrow (2, 3)_{-.04} \rightsquigarrow (3, 3)_{-.04} \rightsquigarrow (3, 2)_{-.04} \rightsquigarrow (3, 3)_{-.04} \rightsquigarrow (4, 3)_{+1}$$

$$(1, 1)_{-.04} \rightsquigarrow (2, 1)_{-.04} \rightsquigarrow (3, 1)_{-.04} \rightsquigarrow (3, 2)_{-.04} \rightsquigarrow (4, 2)_{-1}$$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them



# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them

Agent	We Know	What to Learn	What to use
-------	---------	---------------	-------------

---

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$
Q-learning agent		$Q(s, a)$	$Q$

# Agent Types for RL

## 3 Reinforcement Learning based Agents

What to do if we do not know  $R$ , or  $P$ , or both?

- We can try to learn  $R$  and  $P$
- Or something else that obviates the need to learn them

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$
Q-learning agent		$Q(s, a)$	$Q$
Reflex agent		$\pi(s)$	$\pi$

## Passive vs. Active

### 3 Reinforcement Learning based Agents

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$
Q-learning agent		$Q(s, a)$	$Q$
Reflex agent		$\pi(s)$	$\pi$

- Besides these designs, we decide how “adventurous” our agent will be

## Passive vs. Active

### 3 Reinforcement Learning based Agents

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$
Q-learning agent		$Q(s, a)$	$Q$
Reflex agent		$\pi(s)$	$\pi$

- Besides these designs, we decide how “adventurous” our agent will be

**Passive RL** The agent simply follows a fixed policy, and learns the components above as he goes  
The goal is to learn how good that policy is

## Passive vs. Active

### 3 Reinforcement Learning based Agents

Agent	We Know	What to Learn	What to use
Utility-based agent	$P$	$R \rightarrow U$	$U$
Q-learning agent		$Q(s, a)$	$Q$
Reflex agent		$\pi(s)$	$\pi$

- Besides these designs, we decide how “adventurous” our agent will be

**Passive RL** The agent simply follows a fixed policy, and learns the components above as he goes  
The goal is to learn how good that policy is

**Active RL** After learning the components for a while,  
the agent changes policy to “cash in” on the learned components

To continue in the next session.