

# An Introduction to Reinforcement Learning

Hossein Hajiabolhassan

Department für Mathematik  
und Informationstechnologie  
Montanuniversität Leoben

Zhejiang Normal University Seminar

November 11, 2022



# Table of Contents

- ① Preliminary
- ② Training an Agent: Assigning Rewards
- ③ Challenges and Limitations
- ④ Small Environment and Multi-Armed Bandits
- ⑤ Various Algorithms for Multi-Armed Bandits
- ⑥ Applications and Variants

# Preliminary

# Reinforcement Learning

## Definition (Reinforcement Learning)

Training: Reinforcement Learning is the science of **decision making**. It is about learning the optimal behavior in an environment to obtain **maximum reward and minimum punishment**.

Task: Find a model to maximize the total cumulative reward of the agent.

Examples: Teach an artificial intelligence (AI) system to play GO

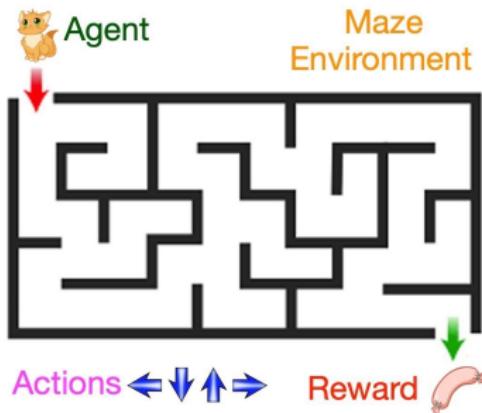


▶ Reference

# Reinforcement Learning

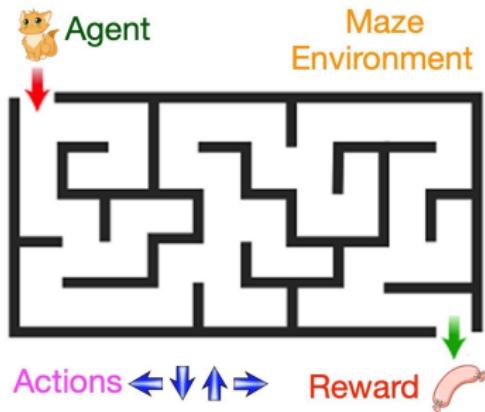
## Training

How is it possible to **train** a **cat** to find **exit** easily?



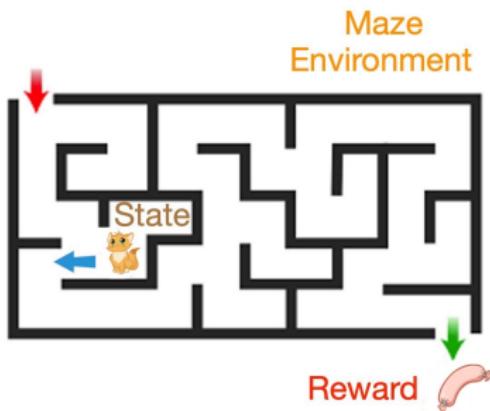
## Reinforcement Learning

- **Agent**: The learner or decision-maker.
  - **Environment**: The Agent's world in which it lives and interacts.
  - The goal of reinforcement learning is to train an agent to complete a task within environment



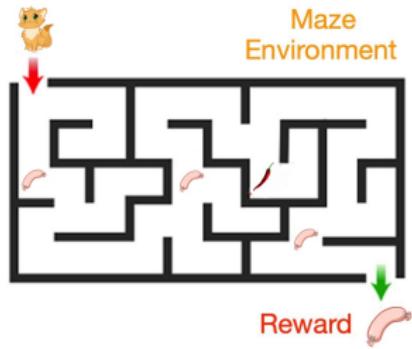
# Reinforcement Learning

- **Action:** The set of all possible moves an agent can make. The agent can interact with the environment by performing some action but cannot influence the rules or dynamics of it by those actions.
- **State:** An immediate situation in which the agent finds itself. It can be a specific moment or position in the environment.



## Reinforcement Learning

- **Reward**: For every action made, the agent receives a reward from the environment and indicates performance of the agent.
  - **Reward Hypothesis**: All goals can be described by the optimization of some function of the sequence of rewards.
  - Agent is not so clever to learn alone. Agent's goal is to optimize some function of the sequence of rewards to reach the goal.



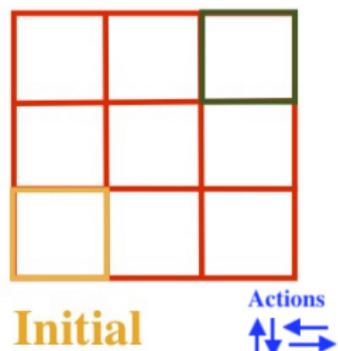
# Training an Agent: Assigning Rewards

# Goal

- **Shortest Path:** Train an agent which is in **the initial state** to reach to the **terminal state** in the shortest possible time.
- How is it possible to **train the agent** to reach the goal?

- ① **Terminal State** is an absorbing state. All actions taken in the absorbing state lead back to that same state.
- ② Agent learns by doing actions and **interacting with environment** (collecting data: trial and error, former experiences, ...)
- ③ By assigning **appropriate rewards**, an agent must be able to learn from **its own experiences**.

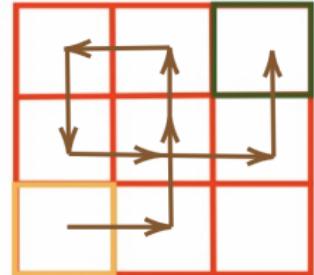
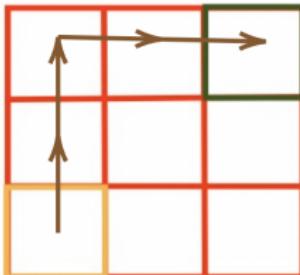
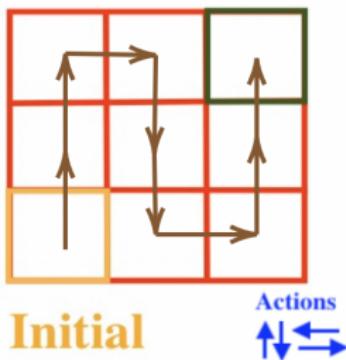
Terminal



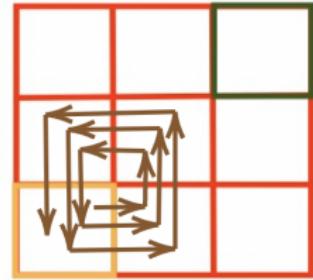
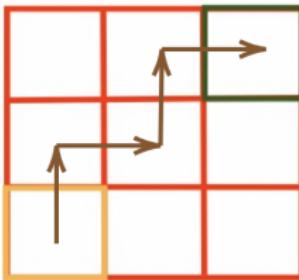
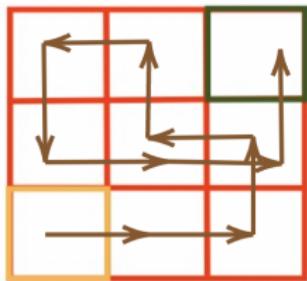
Initial



# Experiences

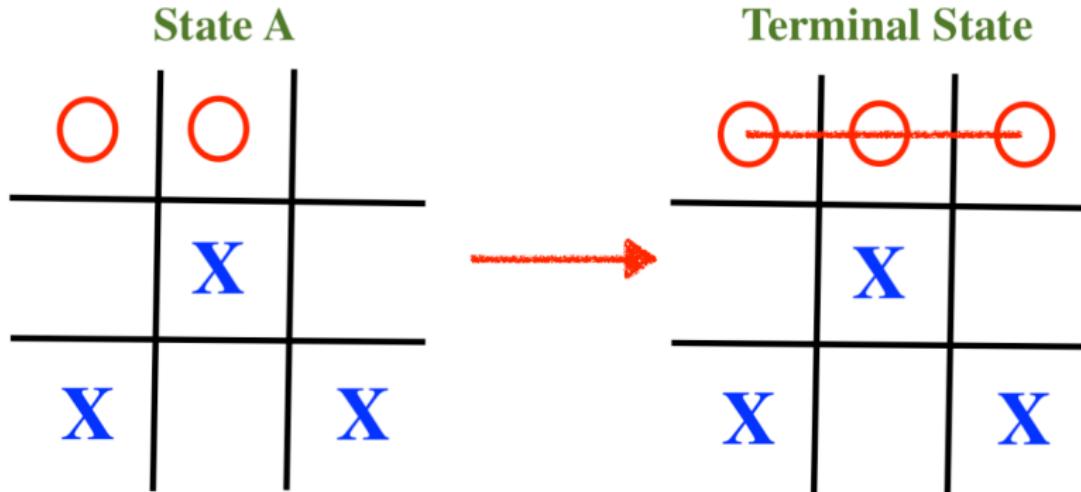


Terminal



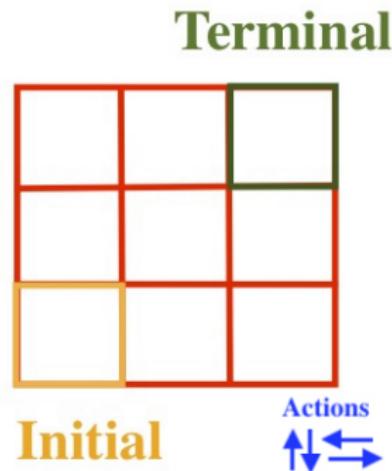
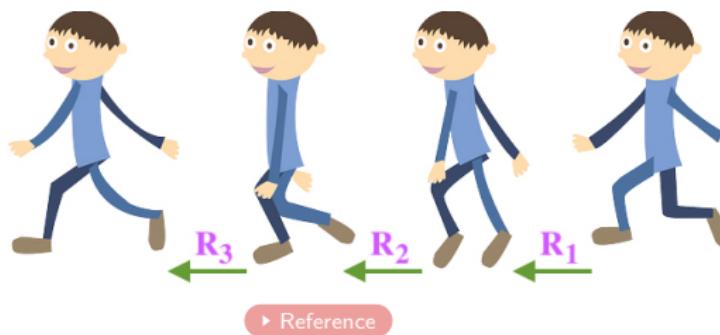
# Experiences

For instance in Tic-Tac-Toe, collecting data does not mean to have all of wining strategies! Also, learning does not mean to copy a wining strategy!



# Choosing a Reward Function

- We usually consider **the cumulative sum of rewards** as a function to reach the goal (Cumulative Sum:  $R_1 + R_2 + \dots + R_n$ ). In this case, the **existence of the terminal state** is necessary.



# Choosing a Reward Function

What is the best candidate?

- Cumulative Sum:  $R_1 + R_2 + \dots + R_n$
- All reward on JUST terminal state (goal)!

## Terminal

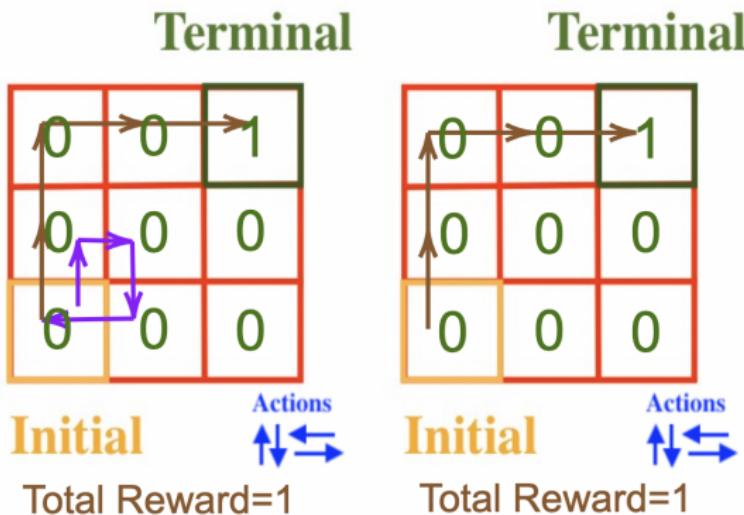
0	0	1
0	0	0
0	0	0

## Initial

# Choosing a Reward Function

What is the best candidate?

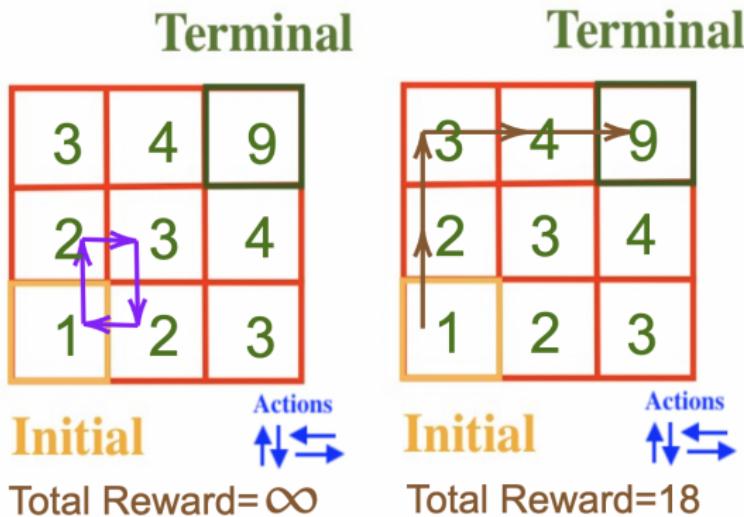
- Cumulative Sum:  $R_1 + R_2 + \dots + R_n$
- All reward on JUST terminal state (goal)!



# Choosing a Reward Function

What is the best candidate?

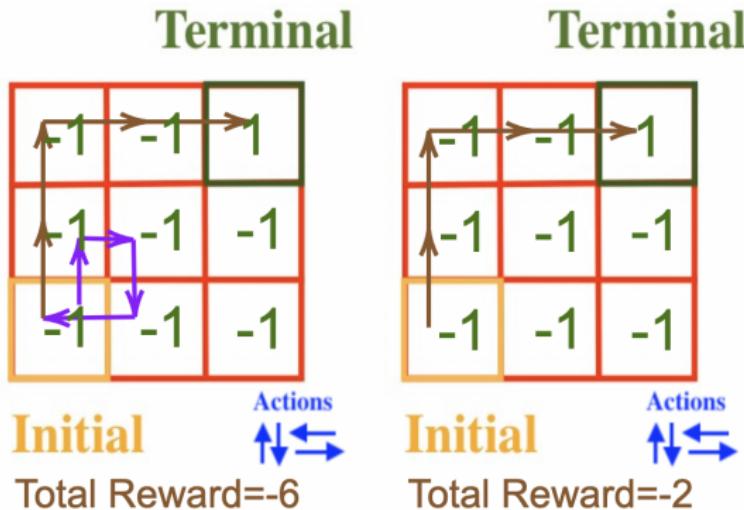
- Increasing reward as you approach terminal states.
- If there exists a cycle, agent can stay in the cycle forever or if the number of steps is large enough, then the summation will be a very large number!



# Choosing a Reward Function

What is the best candidate?

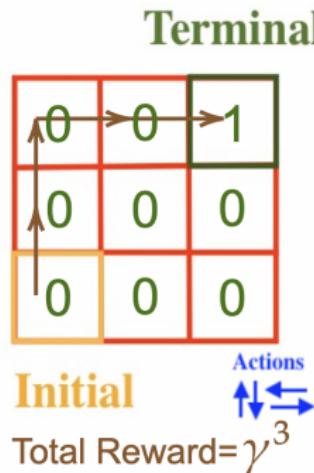
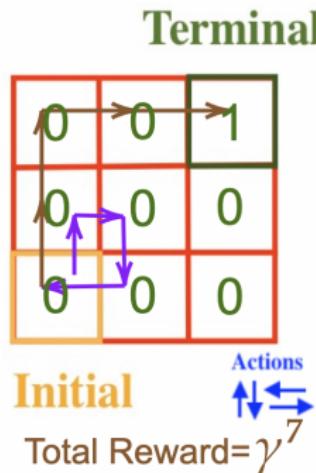
- What is the **solution**?
- Negative reward everywhere (except actions of terminal state)



# Choosing a Reward Function

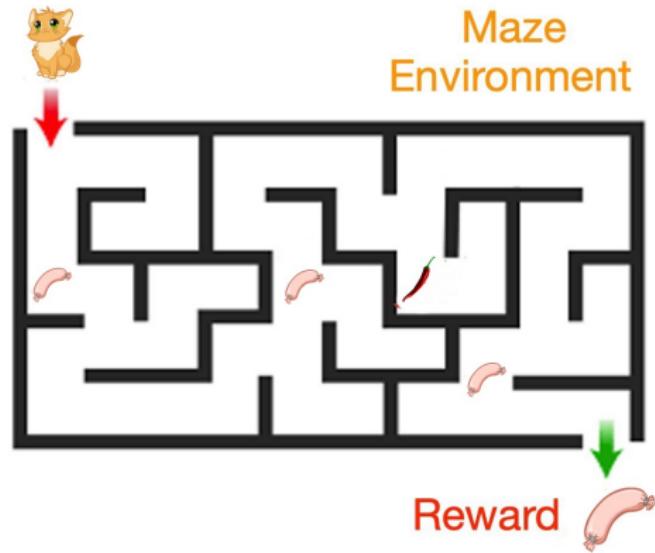
What is the best candidate?

- Valuing immediate money more than future money is a rational behavior known as **discounting**.
- We can consider the **discounting rewards**  $\sum_{i=1}^t \gamma^{i-1} R_i$  where  $0 < \gamma < 1$ .



## Algorithm

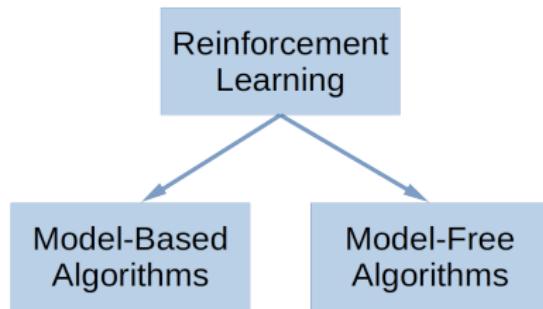
- We assign the rewards such that it is possible to **train** an agent. In training an **AI**, we should consider an appropriate **algorithm**.



# Computational Strategies: Challenges and Limitations

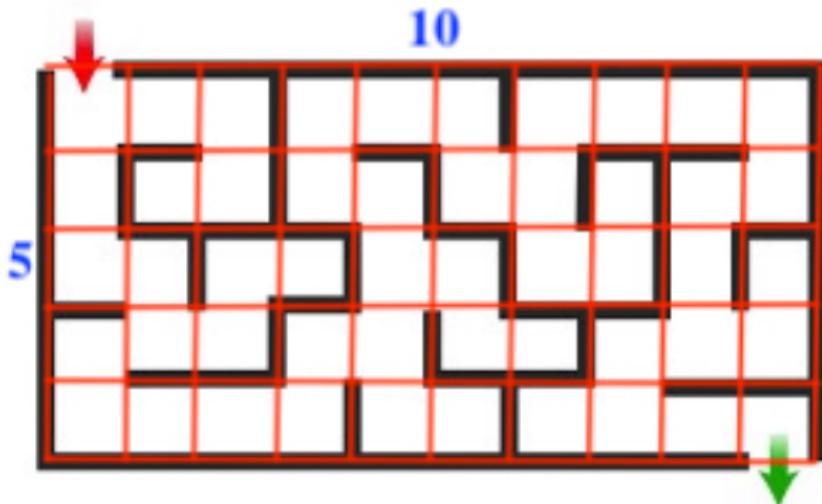
## Models

- ① **Model-based algorithms** use a predictive model of the environment to ask questions of the form “what will happen if I do action  $a$ ?” to choose the best action  $a$
- ② **Model-free algorithms** seek to learn the consequences of their actions through experience



# Known Environment

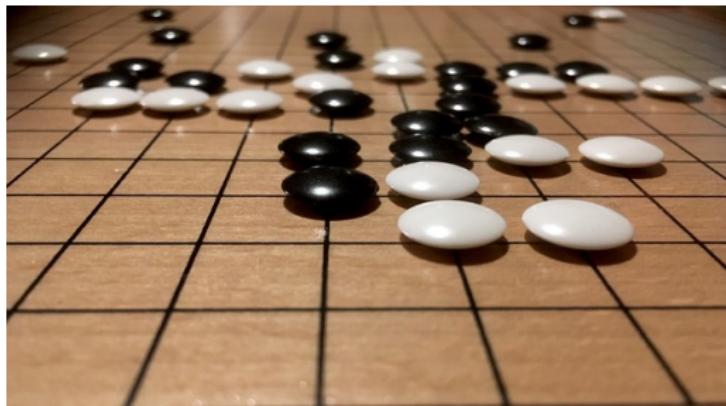
- **Question:** How many states has this environment?
- **Challenge:** The environment is usually **large**.



► Reference

# Known Environment

- ① In GO, the number of **possible board positions** is around  $10^{172}$ .
- ② There are between  $10^{78}$  to  $10^{82}$  atoms in the observable **universe**.



► Reference

# Unknown Environment

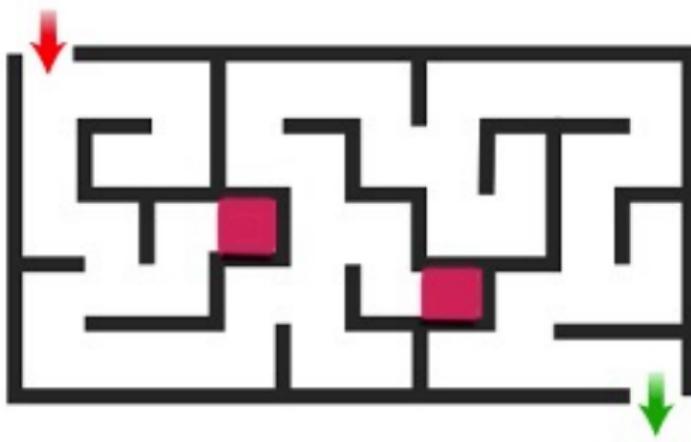
- At the beginning of exploring of the environment, if the agent do an action, then the agent moves to an **unknown state**.
- For instance, when we are **playing a game** and we should choose **a door** to continue the game and we are not aware about next state.



▶ Reference

# Unknown Environment

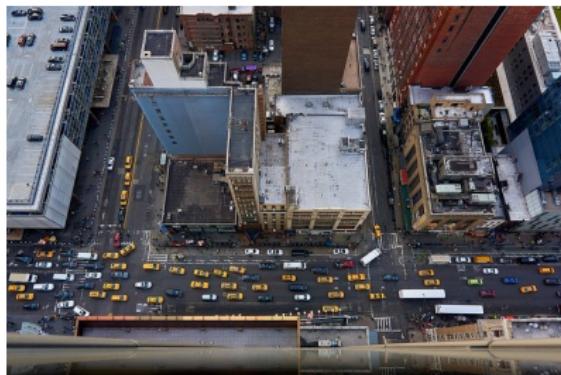
- Is it possible for the agent to **distinguish** between the **red states**? If we assume that all of walls are the same and high enough respect to the agent (the agent does not see the objects behind the wall).



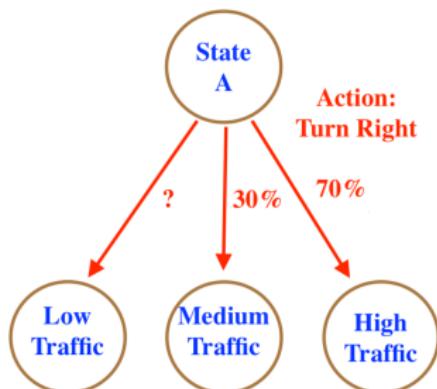
► Reference

# Unknown Environment

- Actions may be stochastic and one may be unlucky to explore **all of transitions**. For instance, in the following Environment and in the state **A**, the action **Turn Right** is **stochastic** and the state **Low Traffic** has not been explored as its neighbour so far!

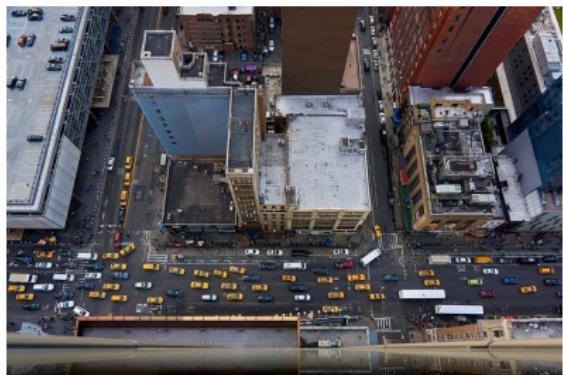


▶ Reference

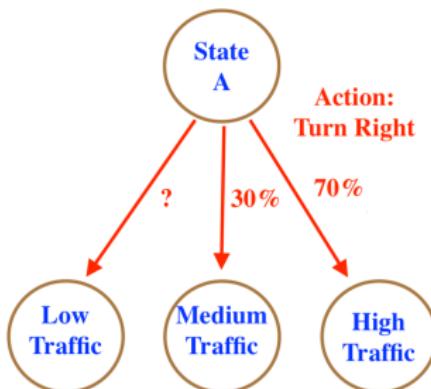


# Unknown Environment

- **Limitation:** The agent usually has no chance to run several actions sequentially in a specific state to **explore** the environment.
- **Limitation:** Agent usually cannot restart to **explore** the environment from **an arbitrary state**.



▶ Reference



# Full Observability

- Challenge: Partial Observability

In this talk, we assume full observability: the new state resulting from executing an action will be known to the system.

An example for partial observability: in some card games (e.g. Bridge), the agent cannot see the hands of the opponent for the most part of the game.



► References

# The Smallest Environment: Multi-Armed Bandits

# Multi-Armed Bandit

**Question:** Is there any challenge for small environments?

# Multi-Armed Bandit

**Question:** Is there any challenge for small environments?

“Bandit” is someone who robs people, especially one of a group of people who attack travellers.



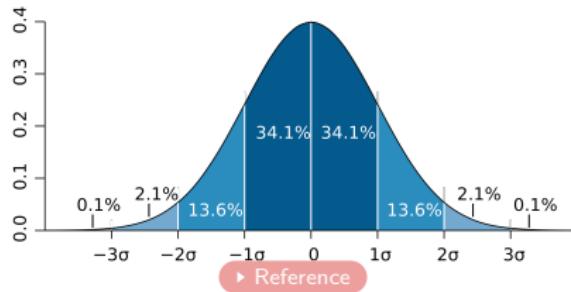
► Reference

“Bandit” in “Multi-Armed Bandits” comes from “one-armed bandit” machines used in a **casino**.



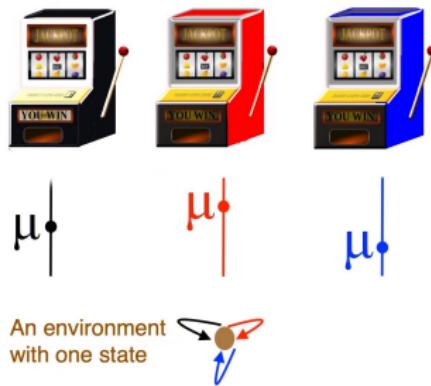
► Reference

- **Question:** Does a teacher always give the same reward when a student solve a kind of problem?
- The rewards and punishments are **often non-deterministic**, and there are invariably stochastic elements governing the underlying situation.
- **Question:** If the reward has **normal distribution** with **unknown mean**, how do we find a approximation of the mean of distribution?
- **The law of large numbers** states that as a sample size grows, its mean gets closer to the average of the whole population.



# Small Environment

- Given a bandits with  $K$  arms where each arm has a fixed but unknown probability distributions with mean  $\mu$ . Pulling any arm, gives you a stochastic reward. In this talk, we usually consider Gaussian Bandits.
- Gaussian Bandits: The rewards come from a normal distribution.
- Bernoulli Bandits: At each round, we receive a binary reward (0 – 1).



▶ Reference

# Goals in Multi-Armed Bandit

- Our Goal can be one of the following items:

- ① Pull the arms one-by-one in sequence such that we **maximize our total reward** collected in the long run.
- ② Best arm identification: **minimize the error probability** at time T.
- ③ Best arm identification: **minimize the total number of stages** used to return the best arm with probability  $1 - \delta$ .

## Training to Reach the Goal

One can use the **rewards** of bandits to **train the agent**.



▶ Reference

# Gaussian Multi-Armed Bandit



Pull arm 2  
0.3



Pull arm 3  
0.1

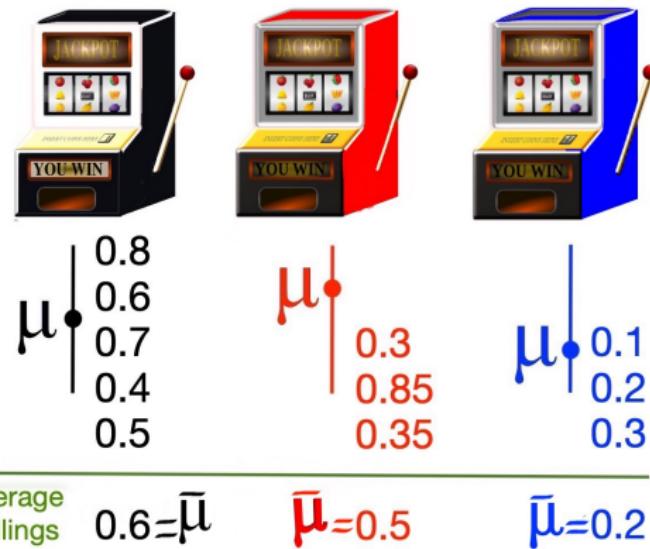


Pull arm 1  
0.8

▶ Reference

# Gaussian Multi-Armed Bandit

- The multi-armed bandit problem



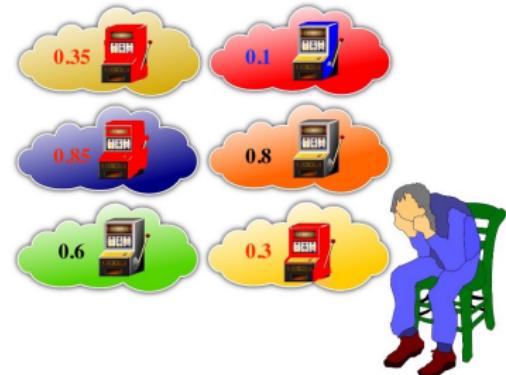
► Reference

# Various Algorithms for Multi-Armed Bandits

# Algorithm and Policy

The agent needs an algorithm to reach his/her goal!

- Goal: maximize the cumulative sum of rewards



## Greedy Algorithm: Based on Exploitation

- 1 Pull all arms once! (Exploration)
  - 2 Choose the best arm!
  - 3 Pull the best arm again (Exploitation)
  - 4 Recompute the average reward of arms
  - 5 Go to step 2 and continue
- What is the best algorithm?

► References 1,

► 2,

► 3

# Performance Measure of Algorithms

- ① There are  $K$ -arms.
- ②  $\mu^* = \max_{1 \leq i \leq K} \mu_i$  denotes the optimal mean of arms.
- ③ In each time step  $t$ , the agent has to play an arm  $l_t \in \{1, 2, \dots, K\}$ .
- ④ Let  $x_{l_t}(t)$  be the obtained reward at time  $t$  where  $1 \leq t \leq T$ .
- ⑤ The aim is to maximize the cumulative sum of rewards.
- ⑥ The rewards are stochastic, so we consider the expectation of them

$$\sum_{t=1}^T E(x_{l_t}(t)) = \sum_{t=1}^T \mu_{l_t}$$

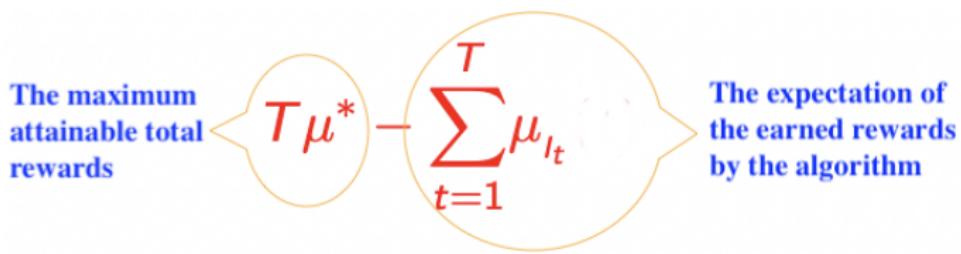
- ⑦ Instead, one can minimize the pseudo regret

$$R(T) = T\mu^* - \sum_{t=1}^T \mu_{l_t}$$

# Performance Measure of Algorithms

- ① There are  $K$ -arms.
- ②  $\mu^* = \max_{1 \leq i \leq K} \mu_i$  denotes the optimal mean of arms.
- ③ In each time step  $t$ , the agent has to play an arm  $l_t \in \{1, 2, \dots, K\}$ .
- ④ Let  $x_{l_t}(t)$  be the obtained reward at time  $t$  where  $1 \leq t \leq T$ .
- ⑤ The aim is to maximize the cumulative sum of rewards.
- ⑥ The rewards are stochastic, so we consider the expectation of them

$$\sum_{t=1}^T E(x_{l_t}(t)) = \sum_{t=1}^T \mu_{l_t}$$



## Greedy Algorithm: Explore-First with parameter $N$

- ① Pull each of  $K$  arms  $N$  times! (Exploration)
- ② Choose the best arm!
- ③ Pull the best arm again (Exploitation)
- ④ Recompute the average reward of arms
- ⑤ Go to step 2 and continue

### Theorem

If  $T$  is known, Explore-first achieves regret:

$$R(T) \leq T^{\frac{2}{3}} \times O(K \log T)^{\frac{1}{3}}.$$

The performance of the exploration phase is terrible.

# Lower Bound for Regret

Theorem (P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire)

*Fix time horizon  $T$  and the number of arms  $K$ . For any bandit algorithm, there exists a problem instance such that*

$$R(T) \geq \Omega(\sqrt{KT}).$$

## Theorem

*If  $T$  is known, Explore-first achieves regret:*

$$R(T) \leq T^{\frac{2}{3}} \times O(K \log T)^{\frac{1}{3}}.$$

The **performance** of the exploration phase is **terrible**.

# The Exploration-Exploitation Dilemma

Challenge

## The Exploration-Exploitation Dilemma

**Exploitation:** make the best decision given current information

**Exploration:** gather more information

- ① The best long-term strategy may involve short- term sacrifices
- ② Gather enough information to make the best overall decision



Reference

# Food Selection in Food Court

**Exploit:** Eat your favourite food



► Reference

**Explore:** Try a new food

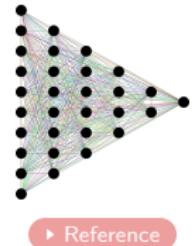


► Reference

# Examples of The Exploration-Exploitation Dilemma

## ① Hyperparameter Tuning in Machine Learning:

**Exploit:** Use the best known Hyperparameters  
**Explore:** Check new Hyperparameters

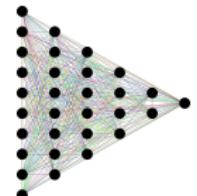


## Examples of The Exploration-Exploitation Dilemma

- ## ① Hyperparameter Tuning in Machine Learning:

**Exploit:** Use the best known Hyperparameters

**Explore:** Check new Hyperparameters



## ► Reference

- ## ② Skill Improvement:

**Exploit:** Use current skills

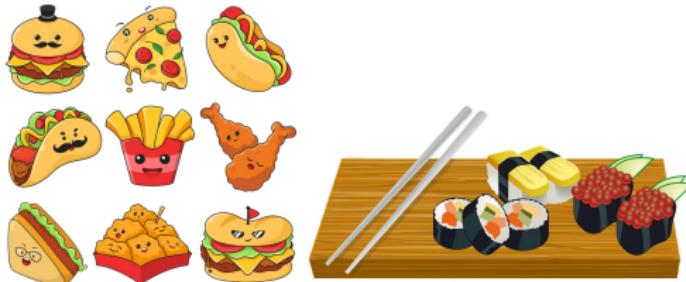
## Explore: Learn new skills



# Different Algorithms: $\epsilon_T$ -Greedy

- **Explore** with probability  $\epsilon_T = T^{-\frac{1}{3}}(K \log T)^{\frac{1}{3}}$
  - **Exploit** with probability  $1 - \epsilon_T$
- ① Pull each of  $K$  arms once!
  - ② At time step  $T$ , pull **the best arm** with probability  $1 - \epsilon_T$ , otherwise pull an **arbitrary arm** with probability  $\epsilon_T$ .
  - ③ Recompute the average reward of arms and Go to step 2.

**Question:** In food selection, we usually use  $\epsilon$ -Greedy algorithm. Which  $\epsilon_T$  do you use?



▶ Reference

▶ Reference

# Different Algorithms: $\epsilon_T$ -Greedy

- **Explore** with probability  $\epsilon_T = T^{-\frac{1}{3}}(K \log T)^{\frac{1}{3}}$
  - **Exploit** with probability  $1 - \epsilon_T$
- ① Pull each of  $K$  arms once!
  - ② At time step  $T$ , pull **the best arm** with **probability**  $1 - \epsilon_T$ , otherwise pull an **arbitrary arm** with **probability**  $\epsilon_T$ .
  - ③ Recompute the average reward of arms and Go to step 2.

## Theorem

*Algorithm  $\epsilon_T$ -Greedy achieves regret:*

$$R(T) \leq T^{\frac{2}{3}} \times O(K \log T)^{\frac{1}{3}}.$$

The **performance** of the exploration phase is still **terrible**.

# Advanced Algorithms (Adaptive Exploration)

- ➊ Pull each of  $K$  arms once!
- ➋ Pull the arm whose empirical average reward prior to time  $t$  plus its error has the maximum value and repeat this step!

$$\max_{\text{arms}} \left( \bar{\mu}_i(t) + \text{Error}_i(t) \right)$$

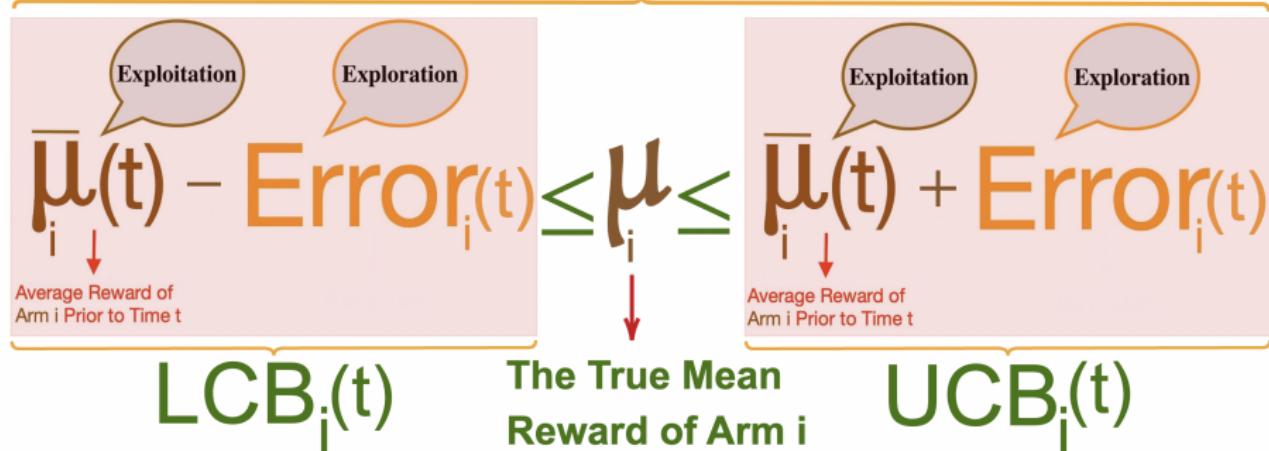
↑  
Average Reward of  
Arm  $i$  Prior to Time  $t$

↑  
 $\beta_i(t)$

The equation represents the selection rule for an arm  $i$  at time  $t$ . It combines the average reward  $\bar{\mu}_i(t)$  (labeled as 'Exploitation') with an exploration term  $\text{Error}_i(t)$  (labeled as 'Exploration'). The exploration term is proportional to the uncertainty or variance of the reward, represented by  $\beta_i(t)$ .

# Advanced Algorithms (Adaptive Exploration)

Concentration Lemma: Hoeffding's Lemma, ...



# Upper Confidence Bound Bandit (UCB1)

- ➊ Pull each of  $K$  arms once!
- ➋ Pull the arm whose empirical average reward prior to time  $t$  plus its error has the maximum value and repeat this step!

$$\max_{\text{arms}} \left( \bar{\mu}_i(t) + \sqrt{\frac{c \log(t)}{\# \text{Pulling of Arm } i}} \right)$$

Average Reward of  
Arm  $i$  Prior to Time  $t$

Exploitation

Exploration

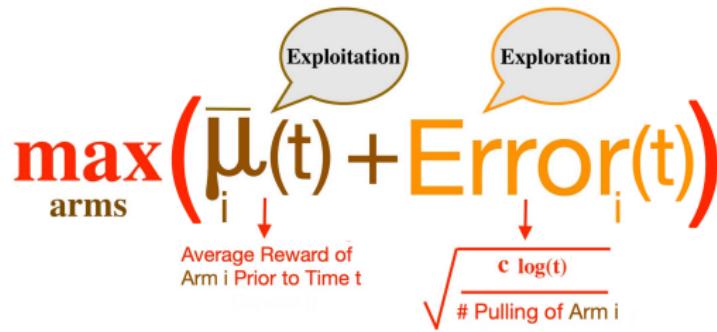
# Upper Confidence Bound Bandit (UCB1)

- ① Pull each of  $K$  arms once!
- ② Pull the arm whose **empirical average reward** prior to time  $t$  plus its **error** has the maximum value and repeat this step!

Theorem (P. Auer, N. Cesa-Bianchi, and P. Fischer)

for all rounds  $t \leq T$ , Algorithm achieves regret:

$$R(T) = O(\sqrt{Kt \log T}).$$



# Bayesian (Based Error)

- ➊ Pull each of  $K$  arms once!
- ➋ Pull the arm whose empirical average reward prior to time  $t$  plus its error has the maximum value and repeat this step!

$$\max_{\text{arms}} \left( \bar{\mu}_i(t) + \text{Error}_i(t) \right)$$

↑  
Exploitation      Exploration  
↓  
Average Reward of Arm  $i$  Prior to Time  $t$       Bayesian

# Bayesian Algorithms

- ① Pull each of  $K$  arms once!
- ② Pull the arm whose empirical average reward prior to time  $t$  plus its error has the maximum value and repeat this step!
- ③ Bayesian:  $\max_i \{\bar{\mu}_i(t) + \bar{\sigma}_i(t) \times \Phi^{-1}(1 - \alpha_t)\}$ , where  $\bar{\mu}_i(t)$  and  $\bar{\sigma}_i(t)$  are the Bayesian mean, Bayesian variance. Also,  $\Phi^{-1}(\cdot)$  stands for quantile function of the standard normal random variable!

$$\max_{\text{arms}} \left( \bar{\mu}_i(t) + \text{Error}_i(t) \right)$$

↓  
Average Reward of  
Arm  $i$  Prior to Time  $t$

↓  
Bayesian



# Finite Regret

Question: is it possible to introduce an algorithm with finite regret? NO!

Theorem (P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire)

Fix time horizon  $T$  and the number of arms  $K$ . For any bandit algorithm, there exists a problem instance such that

$$R(T) \geq \Omega(\sqrt{KT}).$$

# Finite Regret

Question: is it possible to introduce an algorithm with finite regret? NO!

- ① In practice, we are often happy to perform a task just good enough.  
For example, when driving to work we will be content with a strategy that will let us arrive just in time.
- ② We only care about whether an arm with mean reward  $\geq S$  is chosen, where  $S$  is the level of satisfaction we aim at.
- ③ Modify the classic notion of regret and consider, the satisfying (pseudo-)regret with respect to  $S$  (short  $S$ -regret) defined as

S-Regret:

$$R_S(T) = \sum_{t=1}^T \max\{S - \mu_{I_t}, 0\}$$

Regret:

$$R(T) = T\mu^* - \sum_{t=1}^T \mu_{I_t}$$

# Satisfying Bandit

Theorem (T. Michel, H.H., R. Ortner)

If  $S < \mu^*$ , there exists a constant  $C$  such that for any  $T$ ,  $R_S(T) \leq C$ .

## ALGORITHM

**Input:**  $K, S$

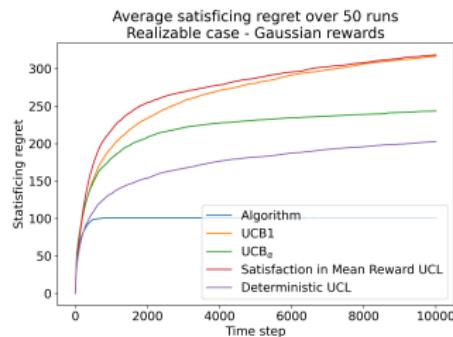
- ▷ Play each arm once.
- ▷ **For** each further step  $t$  **do**

**If**  $\exists$  arm  $i$  with  $\hat{\mu}_i \geq S$  play

$$A_t \in \operatorname{argmax}_{1 \leq i \leq k} \left\{ \frac{\text{UCB}_i(t) - \max\{S, \text{LCB}_i(t)\}}{\beta_i(t)} \right\}.$$

**Else If**  $\exists$  arm  $i$  with  $\text{UCB}_i(t) \geq S$  choose  $A_t$  uniformly at random from  $\{i | \text{UCB}_i \geq S\}$ .

**Else** choose  $A_t \in \operatorname{argmax}_{1 \leq i \leq k} \text{UCB}_i(t)$ .



▶ Reference

# Best Arm Identification: Two Arms

## Exploration versus Exploitation

Based on the goal, the optimal algorithms may move towards **more exploration** than **exploitation**

- Assume that there are just **two arms**.
- The distribution of the rewards of each arm is Gaussian with unknown mean and **the same variance**.
- The goal is the **best arm identification**: **minimize the error probability** at **time  $2T$** .

### Algorithm:

- ① Pull each arm  **$T$  times!**

# Applications and Variants

# Applications

- ① K treatments for a given symptom (with unknown effect)
- ② What treatment should be allocated to the next patient, based on responses observed on previous patients?



► Reference

# Applications

- ① K treatments for a given symptom (with unknown effect)
- ② What treatment should be allocated to the next patient, based on responses observed on previous patients?



▶ Reference

- ③ K adds that can be displayed
- ④ Which add should be displayed for a user, based on the previous clicks of previous (similar) users?



▶ Reference



# Contextual Bandit

- ① Assign to each person a **feature** or a **context**.
- ② Imagine that each **machine** responds differently to **each person**.
- ③ You need to find **the best strategy** for the given **context**.
- ④ **K treatments** and **K ads** can be considered as **contextual bandit**.
- ⑤ Variants: **Adversarial Bandit**, **Infinite-Armed Bandit**, **Non-Stationary Bandit**, **and so on**



▶ Reference

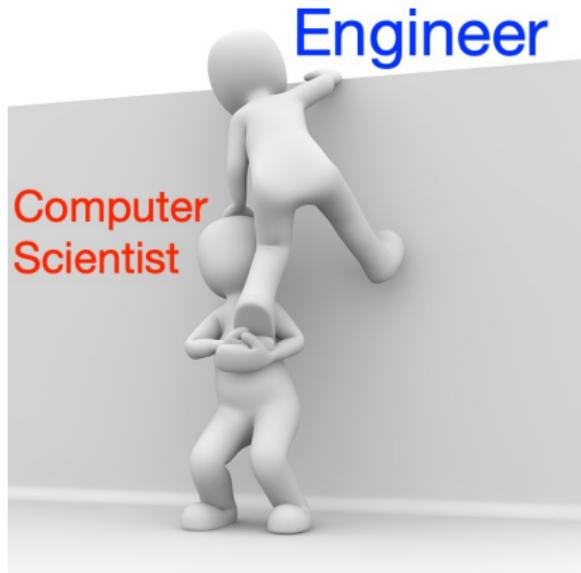


▶ Reference

**Ultimate Goal:** For given problem, find algorithm that works well in practice (performance, efficiency), has favorable theoretical guarantees.

**Engineer:** Your theoretical results helped me so much about  $K$  adds!

**Computer Scientist:** You are welcome!



▶ Reference

# Acknowledgement

Thanks for your attention!

I would like to show my greatest appreciation to Professor Ronald Ortner and Professor Xuding Zhu for the immeasurable amount of their support. Furthermore, they challenged me to push further and to learn more.

