

**LSE** **ST455: Reinforcement Learning**  
**Lecture 1: Introduction to Reinforcement Learning**

Chengchun Shi

# Lecture Outline

---

1. Introduction and Course Overview
2. Multi-Armed Bandit
3. Contextual Bandits

# 1. Introduction and Course Overview

## 2. Multi-Armed Bandit

## 3. Contextual Bandits

# Course Information

---

- **Lectures:** Tue 14 – 16:00 pm, PAR 1.02 (zoom: 985 785 4435)
- **Seminars:** Wed 13 – 14:30 pm, CBG 1.05 (zoom: 985 785 4435; lead by CS)  
Fri 15 – 16:30 pm, CLM 5.02 (lead by Domenico)

You may join lecture and seminar via **zoom** as well

- **Office Hours:**

- Chengchun Shi (c.shi7@lse.ac.uk): Tue & Wed 10-11:00 am, COL 8.08 or ZOOM
- Domenico Mergoni (d.mergoni@lse.ac.uk): Fri 10:30-11:30 am
- Please use **LSE Student Hub** to book slots

- **Assessment:**

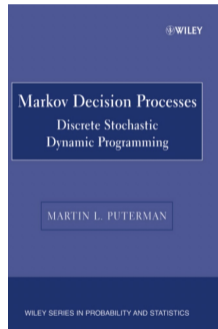
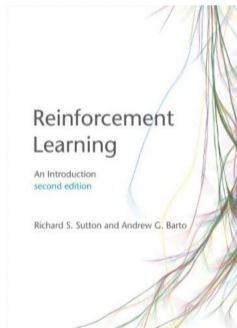
- Two summative assignments at Weeks 4 & 7 (10% each)
- A final project (group project) to apply/develop RL algorithms (80%)

- We use **GitHub**. Please register and fill in the form
- More on **Moodle** ([link](#))

# Textbooks

---

- **Reinforcement Learning: An Introduction** (Second Edition) by Sutton and Barto (2018)
  - **Hardcover** £50 on Amazon
  - **Ebook** free online ([link](#))
  - 50K citations so far
- **Markov decision processes: discrete stochastic dynamic programming** by Puterman (2014)



# Useful Resources

---

- **Deepmind & UCL** reinforcement learning (RL) course by David Silver
  - **Course webpage** [link](#)
  - **Videos** available on Youtube
  - **Slides** available on webpage
- **UC Berkeley** PhD-level deep RL course by Sergey Levine
  - **Course webpage** [link](#)
  - Some more **resources** [link](#)
- Working draft on “**Reinforcement Learning: Theory and Algorithms**” by Alekh, Nan, Sham and Wen [link](#)



# Applications

---



(a) Games



(b) Health Care



(c) Ridesharing



(d) Robotics



(e) Finance



(f) Automated Driving

# Games

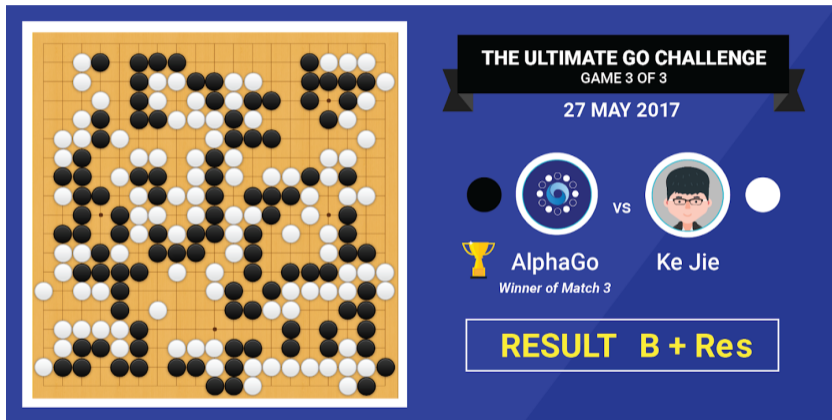


Figure: AlphaGo. See Silver et al. [2016] for details. To be discussed in more detail in Lecture 9.



## Games (Cont'd)

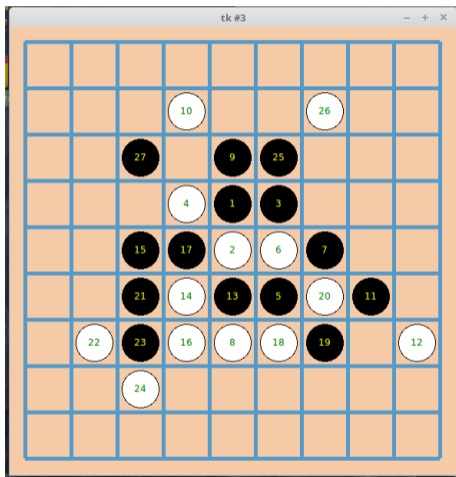


Figure: An implementation of AlphaGo Zero on Gomoku. To be discussed in more detail in Seminar 10.

# Games (Cont'd)

---

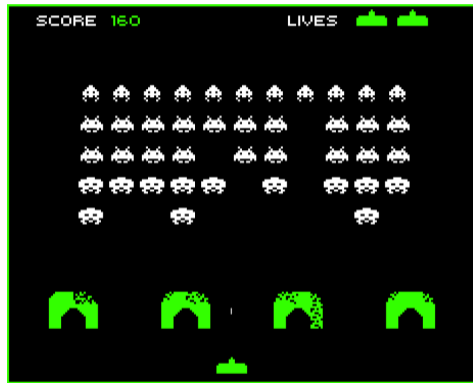
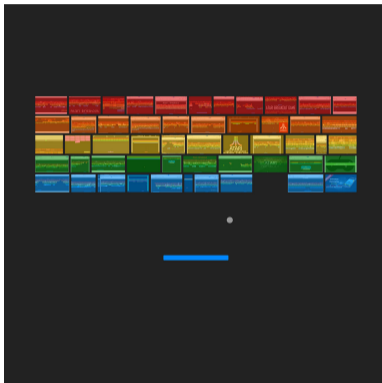


Figure: Two Atari Games: Breakout ([link](#)) and Space Invaders. To be discussed in more detail in Lecture 7 & Seminar 8.

# Healthcare

- Management of **Type-I diabetes** [Luckett et al., 2019, Shi et al., 2020, 2022, Zhou et al., 2022a]
- **Subject:** Patients with Type-I diabetes
- **Objective:** Improve health outcomes
- **Intervention:** Determine whether a patient needs to **inject insulin or not** based on their glucose levels, food intake, exercise intensity, etc.
- **Data:** OhioT1DM dataset

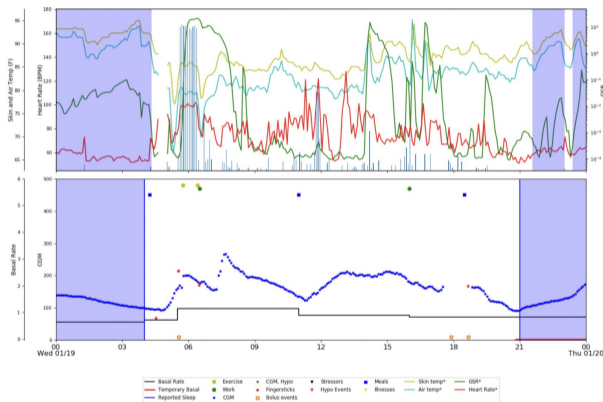


Figure: OhioT1DM data. To be discussed in Lecture 10.

# Healthcare (Cont'd)

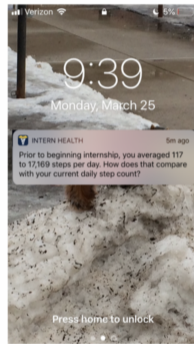
- **Intern health study** [NeCamp et al., 2020, Li et al., 2022]
- **Subject:** First-year medical interns working in stressful environments (e.g., long work hours and sleep deprivation)
- **Objective:** Promote physical and mental well-beings
- **Intervention:** Determine whether to send certain text message to a subject



(i) App Dashboard



(ii) Mood EMA



(iii) Notifications

Figure: IHS. To be discussed in Lecture 10.

# Healthcare (Cont'd)

---

**Table 1.** Examples of 6 different groups of notifications.

Notification groups	Life insight	Tip
Mood	Your mood has ranges from 7 to 9 over the past 2 weeks. The average intern's daily mood goes down by 7.5% after intern year begins.	Treat yourself to your favorite meal. You've earned it!
Activity	Prior to beginning internship, you averaged 117 to 17,169 steps per day. How does that compare with your current daily step count?	Exercising releases endorphins which may improve mood. Staying fit and healthy can help increase your energy level.
Sleep	The average nightly sleep duration for an intern is 6 hours 42 minutes. Your average since starting internship is 7 hours 47 minutes.	Try to get 6 to 8 hours of sleep each night if possible. Notice how even small increases in sleep may help you to function at peak capacity & better manage the stresses of internship.

- Other applications:
  - HeartSteps [Liao et al., 2020]
  - Sepsis treatment [Li et al., 2020, Chen et al., 2022, Zhou et al., 2022b]

# Ridesharing

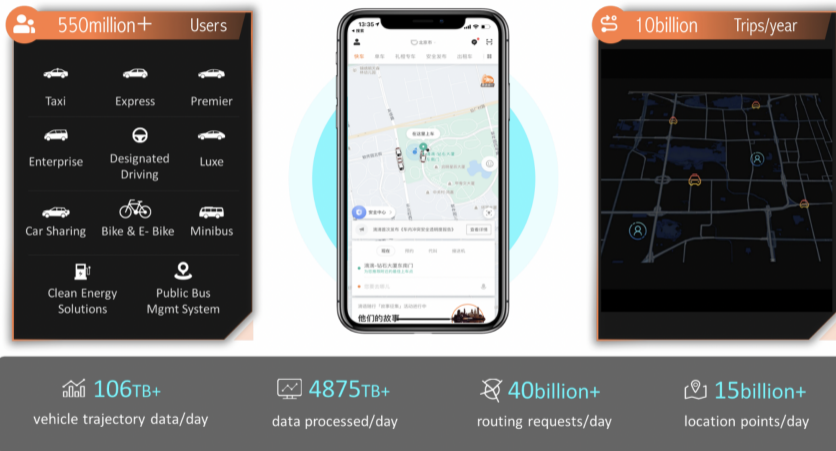
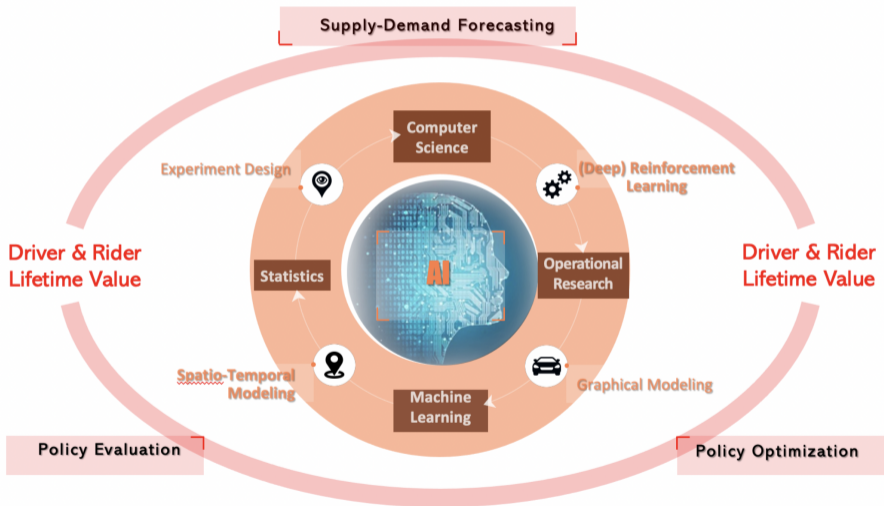


Figure: Ridesharing. To be discussed in more detail in Lecture 7 & Seminar 7.

# Ridesharing (Cont'd)



# Robotics



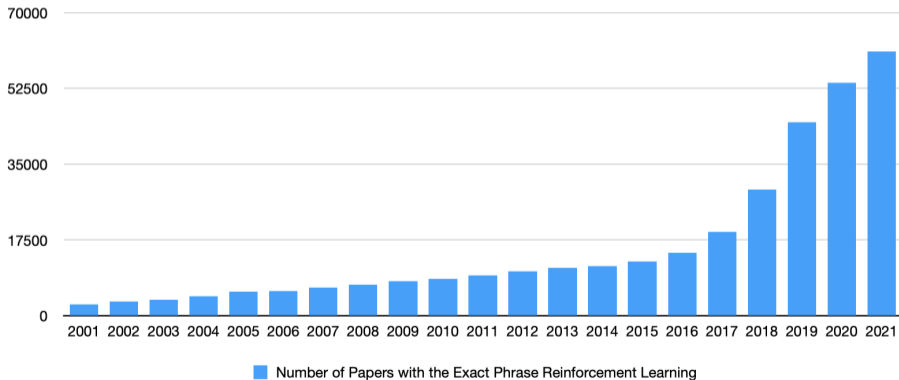
Figure: See <https://www.youtube.com/watch?v=gn4nRCC9TwQ>



# RL as a Research Topic

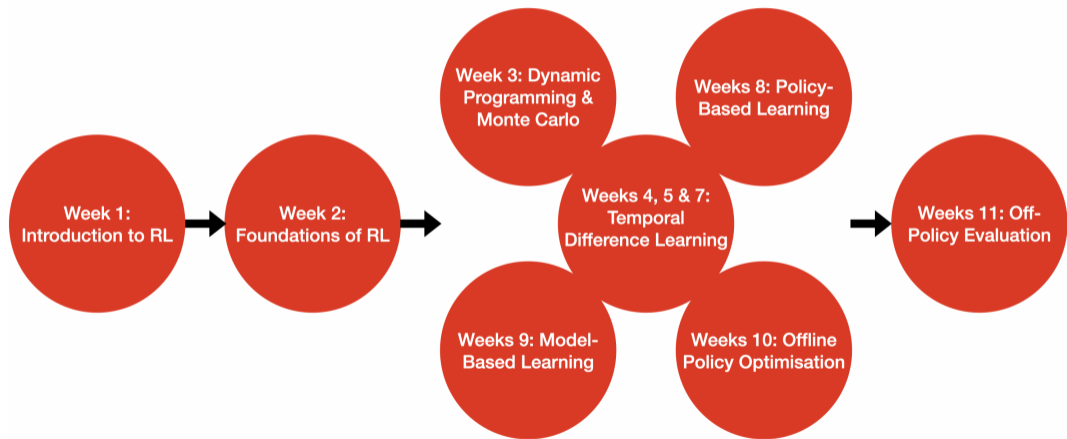
---

- One of the most vibrant research topics in **machine learning**
- Over 100 papers accepted at **ICML 2020**, accounting for more than 10% in total



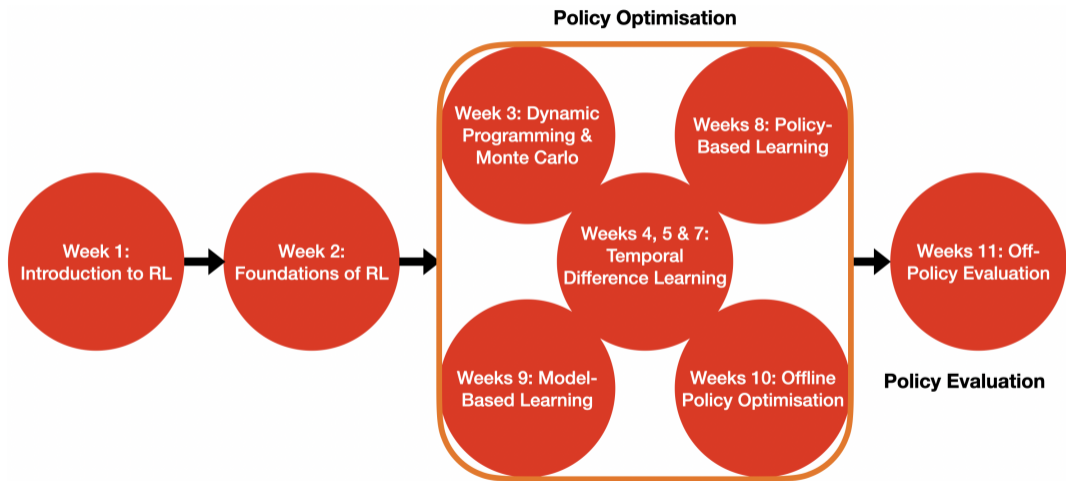
# Roadmap

---



# Roadmap (Cont'd)

---



1. Introduction and Course Overview

**2. Multi-Armed Bandit**

3. Contextual Bandits

# Multi-Armed Bandit (MAB) Problem

---

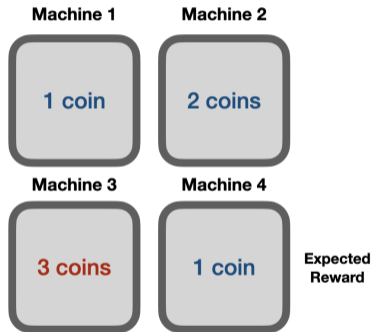


- The **simplest** RL problem
- A casino with **multiple** slot machines
- Playing each machine yields an independent **reward**.
- Limited knowledge (unknown reward distribution for each machine) and resources (**time**)
- **Objective**: determine which machine to pick at each time to maximize the expected **cumulative rewards**

# Multi-Armed Bandit Problem (Cont'd)

---

- $k$ -armed bandit problem ( $k$  machines)
- $A_t \in \{1, \dots, k\}$ : arm (machine) pulled (experimented) at time  $t$
- $R_t \in \mathbb{R}$ : reward at time  $t$
- $Q(a) = \mathbb{E}(R_t | A_t = a)$  expected reward for each arm  $a$  (**unknown**)
- **Objective**: maximize  $\sum_{t=1}^T \mathbb{E}R_t$ .



# Greedy Action Selection

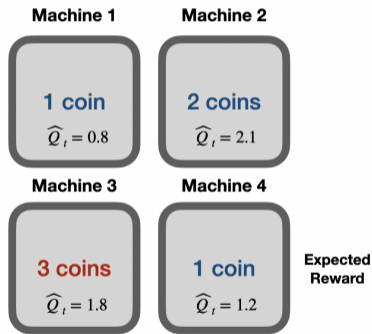
- **Action-value methods:** estimate the expected reward (i.e., value) of actions and use these estimates to select actions
- Estimated reward at time  $t$ :

$$\hat{Q}_t(\mathbf{a}) = \frac{\sum_{i=1}^t R_i \mathbb{I}(\mathbf{A}_i = \mathbf{a})}{\sum_{i=1}^t \mathbb{I}(\mathbf{A}_i = \mathbf{a})}$$

- **Greedy policy:**

$$\mathbf{A}_t = \arg \max_{\mathbf{a}} \hat{Q}_{t-1}(\mathbf{a}).$$

- Might be **suboptimal** in the long run.



# Exploration-Exploitation Dilemma

---

- **Exploitation:** To maximize reward, the agent prefers the greedy policy that selects actions that maximizes the estimated expected reward.
- **Exploration:** To discover which actions yield a higher reward, the agent must try actions that it has less selected to improve the estimation accuracy.
- **Trade-off** between exploration and exploitation:
  - Neither exploration nor exploitation can be used exclusively.
  - The agent must try various actions and progressively favour high-reward actions.
- Practical algorithms:  **$\epsilon$ -greedy**, **upper confidence bound (UCB)**, **Thompson sampling**.



# $\epsilon$ -Greedy

---

- **Input:** Choose a small value parameter  $\epsilon \in (0, 1)$ .
- At each step **perform:**
  - With probability  $1 - \epsilon$ : adopt the **greedy policy**;
  - With probability  $\epsilon$ : choose a **randomly selected arm** from the set of all arms.
- Combines exploration and exploitation:
  - At each time, each arm is selected with probability at least  $k^{-1}\epsilon$ .
  - Greedy action is selected with probability  $1 - \epsilon + k^{-1}\epsilon$ .

# Incremental Implementation

---

- Average reward received from arm  $\mathbf{a}$  by time  $t$ :

$$\hat{Q}_t(\mathbf{a}) = N_t^{-1}(\mathbf{a}) \sum_{i=1}^t \mathbb{I}(\mathbf{A}_i = \mathbf{a}) R_i,$$

where  $N_t(\mathbf{a}) = \sum_{i=1}^t \mathbb{I}(\mathbf{A}_i = \mathbf{a})$ .

- If arm  $\mathbf{a}$  is selected at time  $t + 1$ , then

$$\begin{aligned} \hat{Q}_{t+1}(\mathbf{a}) &= \{N_t(\mathbf{a}) + \mathbf{1}\}^{-1} \left\{ \sum_{i=1}^t \mathbb{I}(\mathbf{A}_i = \mathbf{a}) R_i + R_{t+1} \right\} \\ &= \frac{N_t(\mathbf{a})}{N_t(\mathbf{a}) + \mathbf{1}} \left\{ N_t^{-1}(\mathbf{a}) \sum_{i=1}^t \mathbb{I}(\mathbf{A}_i = \mathbf{a}) R_i \right\} + \frac{R_{t+1}}{N_t(\mathbf{a}) + \mathbf{1}} \\ &= \frac{N_t(\mathbf{a})}{N_t(\mathbf{a}) + \mathbf{1}} \hat{Q}_t(\mathbf{a}) + \frac{R_{t+1}}{N_t(\mathbf{a}) + \mathbf{1}}. \end{aligned}$$

# Algorithm

---

- **Input:**  $0 < \epsilon < 1$ , termination time  $T$ .
- **Initialization:**  $t = 0$ ,  $\hat{Q}(\mathbf{a}) = 0$ ,  $N(\mathbf{a}) = 0$ , for  $\mathbf{a} = 1, 2, \dots, k$ .
- **While**  $t < T$ :
  - **Update**  $t$ :  $t \leftarrow t + 1$ .
  - $\epsilon$ -greedy action selection:

$$\mathbf{a}^* \leftarrow \begin{cases} \arg \max_{\mathbf{a}} \hat{Q}(\mathbf{a}), & \text{with probability } 1 - \epsilon, \\ \text{random arm,} & \text{with probability } \epsilon. \end{cases}$$

- **Receive reward**  $R$  from arm  $\mathbf{a}^*$ .
- **Update**  $N(\mathbf{a}^*)$ :  $N(\mathbf{a}^*) \leftarrow N(\mathbf{a}^*) + 1$ .
- **Update**  $\hat{Q}(\mathbf{a}^*)$ :

$$\hat{Q}(\mathbf{a}^*) \leftarrow \frac{N(\mathbf{a}^*) - 1}{N(\mathbf{a}^*)} \hat{Q}(\mathbf{a}^*) + \frac{1}{N(\mathbf{a}^*)} R.$$

# Example: Four Bernoulli Arms



Bernoulli(0.1)



Bernoulli(0.4)



Bernoulli(0.1)

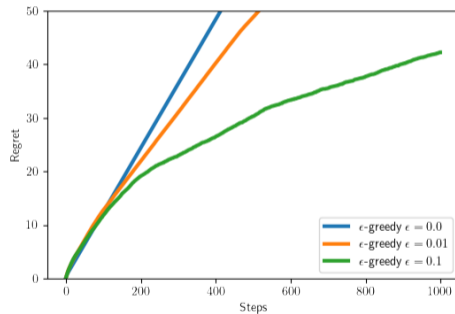
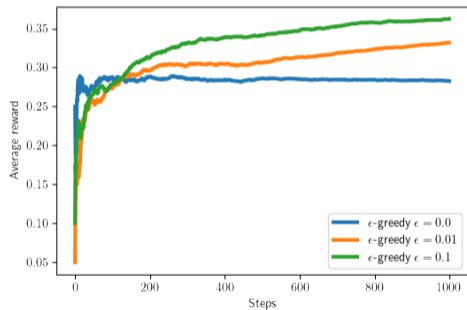


Bernoulli(0.1)

Reward  
distributions

↑  
Best arm

# Example: Four Bernoulli Arms (Cont'd)



# Tracking Nonstationarity

---

- Incremental update:

$$\hat{Q}(\mathbf{a}^*) \leftarrow \frac{N(\mathbf{a}^*) - 1}{N(\mathbf{a}^*)} \hat{Q}(\mathbf{a}^*) + \frac{1}{N(\mathbf{a}^*)} R.$$

- Alternatively, for a given step size parameter  $0 < \alpha < 1$ ,

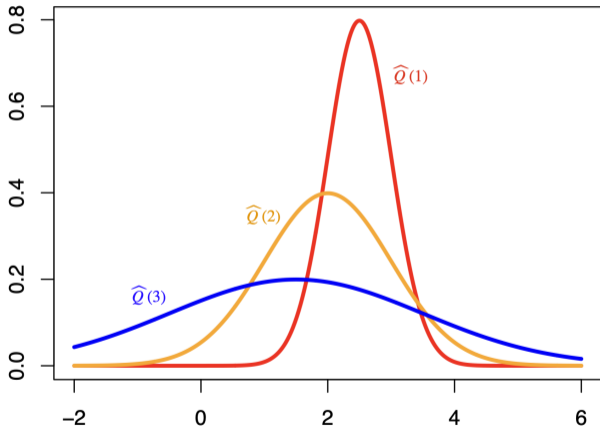
$$\hat{Q}(\mathbf{a}^*) \leftarrow (1 - \alpha) \hat{Q}(\mathbf{a}^*) + \alpha R.$$

- Give more **weight** to recently observed reward. Handles **nonstationarity** (reward distribution varies over time).
- **Exponential weighted moving average:**

$$\begin{aligned} \hat{Q}(\mathbf{a}^*) &\leftarrow \alpha R + (1 - \alpha) \hat{Q}^{(-1)}(\mathbf{a}^*) \leftarrow \alpha R + \alpha(1 - \alpha) R^{(-1)} + (1 - \alpha)^2 \hat{Q}^{(-2)}(\mathbf{a}^*) \\ &\leftarrow \alpha R + \alpha \sum_{i=1}^J (1 - \alpha)^i R^{(-i)}. \end{aligned}$$

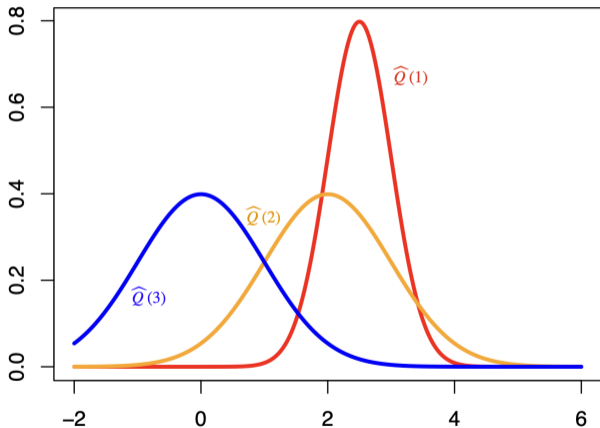
# Optimism in the Face of Uncertainty

- The **optimistic principle**:
- The more **uncertain** we are about an action-value;
- The more **important** it is to explore that action;
- It could be the **best** action.
- Likely to pick blue action.
- **Different** from  $\epsilon$ -greedy which selects arms uniformly random.



# Optimism in the Face of Uncertainty (Cont'd)

- After picking blue action;
- Become less **uncertain** about the value;
- More likely to pick other actions;
- Until we home in on best action.





# Upper Confidence Bound

---

- Estimate an **upper confidence**  $U_t(\mathbf{a})$  for each action value such that

$$Q(\mathbf{a}) \leq \hat{Q}_t(\mathbf{a}) + U_t(\mathbf{a}),$$

with high probability.

- $U_t(\mathbf{a})$  quantifies the **uncertainty** and depends on  $\mathbb{N}_t(\mathbf{a})$  (number of times arm  $\mathbf{a}$  has been selected up to time  $t$ )
  - Large  $\mathbb{N}_t(\mathbf{a}) \rightarrow$  small  $U_t(\mathbf{a})$ ;
  - Small  $\mathbb{N}_t(\mathbf{a}) \rightarrow$  large  $U_t(\mathbf{a})$ .
- Select actions maximizing upper confidence bound

$$\mathbf{a}^* = \arg \max_{\mathbf{a}} [\hat{Q}_t(\mathbf{a}) + U_t(\mathbf{a})].$$

- Combines **exploration** ( $U_t(\mathbf{a})$ ) and **exploitation** ( $\hat{Q}_t(\mathbf{a})$ ).

# Upper Confidence Bound (Cont'd)

---

- Set  $U_t(\mathbf{a}) = \sqrt{c \log(t) / N_t(\mathbf{a})}$  for some positive constant  $c$ .
- According to **Hoeffding's inequality** ([link](#)), when rewards are bounded between  $0$  and  $1$ , the event

$$Q(\mathbf{a}) \leq \hat{Q}_t(\mathbf{a}) + U_t(\mathbf{a}),$$

holds with probability at least  $1 - t^{-2c}$  (converges to  $1$  as  $t \rightarrow \infty$ ).

# Algorithm

---

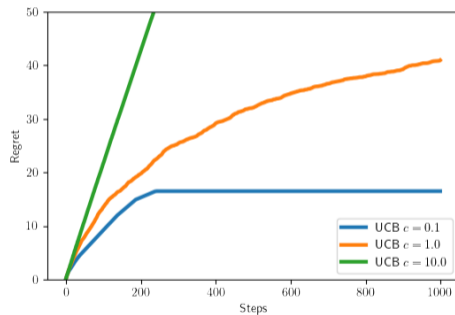
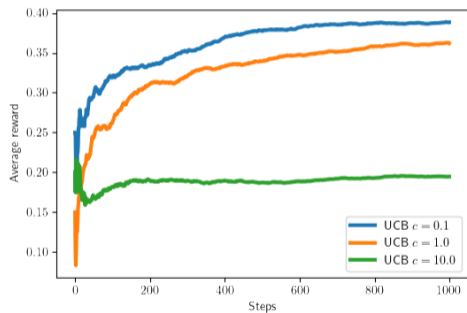
- **Input:** some positive constant  $c$ , termination time  $T$ .
- **Initialization:**  $t = 0$ ,  $\hat{Q}(\mathbf{a}) = 0$ ,  $N(\mathbf{a}) = 0$ , for  $\mathbf{a} = 1, 2, \dots, k$ .
- **While**  $t < T$ :
  - **Update**  $t$ :  $t \leftarrow t + 1$ .
  - **UCB action selection:**

$$\mathbf{a}^* \leftarrow \arg \max_{\mathbf{a}} [\hat{Q}(\mathbf{a}) + \sqrt{c \log(t)/N_t(\mathbf{a})}].$$

- **Receive reward**  $R$  from arm  $\mathbf{a}^*$ .
- **Update**  $N(\mathbf{a}^*)$ :  $N(\mathbf{a}^*) \leftarrow N(\mathbf{a}^*) + 1$ .
- **Update**  $\hat{Q}(\mathbf{a}^*)$ :

$$\hat{Q}(\mathbf{a}^*) \leftarrow \frac{N(\mathbf{a}^*) - 1}{N(\mathbf{a}^*)} \hat{Q}(\mathbf{a}^*) + \frac{1}{N(\mathbf{a}^*)} R.$$

# Example: Four Bernoulli Arms (Revisited)



# Thompson Sampling

---

- A **highly-competitive** algorithm to address exploration-exploitation trade-off.
- Impose **statistical models** for the reward distribution with parameter  $\theta$ .
- Impose **prior distributions** for  $\theta$ .
- At time  $t$ ,
  - Use **Bayes rule** to update the **posterior distribution** of  $\theta$ .
  - Sample a model parameter  $\theta_t$  from the posterior distribution.
  - Compute action-value given  $\theta_t$ , i.e.,  $\mathbb{E}(R|A = a, \theta_t)$ .
  - Select action maximizing action-value

$$a^* = \arg \max_a \mathbb{E}(R|A = a, \theta_t).$$

- Posterior distribution quantifies the **uncertainty** of the estimated model parameter (**exploration**).
- $\mathbb{E}(R|A = a, \theta_t)$  estimates the oracle action value (**exploitation**).

# Thompson Sampling (Cont'd)

---

- **Statistical models:**

- $p(r|\mathbf{a}, \theta)$  models the probability density/mass function of rewards under arm  $\mathbf{a}$ .
- $p(\theta)$  models the probability density/mass function of  $\theta$ .

- **Bayesian inference:**

- Likelihood function  $\ell_t(\theta) = \prod_{i=1}^t p(R_i | A_i, \theta)$ .
- Compute the posterior distribution according to Bayes rule

$$p_t(\theta | \mathcal{D}) = \frac{p(\theta) \ell_t(\theta)}{\int_{\theta} p(\theta) \ell_t(\theta) d\theta} \propto p(\theta) \ell_t(\theta),$$

where  $\mathcal{D}$  denotes the observed data.

- **Compute action value:**

$$\mathbb{E}(R | \mathbf{A} = \mathbf{a}, \theta_t) = \int_r r p(r | \mathbf{a}, \theta_t) dr.$$

# Thompson Sampling (Bernoulli Bandit Example)

---

- **Statistical models:**
  - Reward of the  $a$ th arm follows a Bernoulli distribution with mean  $\theta(a)$ .
  - $\theta(a)$  follows a Beta( $\alpha, \beta$ ) distribution (**prior**).
  - Why Beta distribution?
    - Commonly used distribution for outcomes bounded between **0** and **1**
    - Reduced to **uniform** distribution when  $\alpha = \beta = 1$
    - **Conjugate** distribution of binomial, i.e. posterior distribution is Beta as well
    - $\alpha$  and  $\beta$  measures the beliefs for **success** and **failure**
- **Bayesian inference:**
  - $\theta(a)$  follows a Beta( $S_a + \alpha, F_a + \beta$ ) distribution (**posterior**) where  $(S_a, F_a)$  corresponds to the success and failure counters under arm  $a$ .
- **Compute action value:**

$$\mathbb{E}(R|A = a, \theta_t) = \theta_t(a).$$

# Algorithm (Bernoulli Bandit Example<sup>1</sup>)

---

- **Input:** hyper-parameters  $\alpha, \beta > 0$ , termination time  $T$ .
- **Initialization:**  $t = 0$ ,  $S_a = F_a = 0$ , for  $a = 1, 2, \dots, k$ .
- **While**  $t < T$ :
  - **Update**  $t$ :  $t \leftarrow t + 1$ .
  - **Posterior sampling:** For  $a = 1, 2, \dots, k$ , sample

$$\theta_a \sim \text{Beta}(S_a + \alpha, F_a + \beta)$$

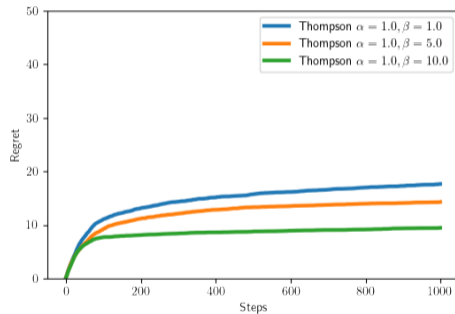
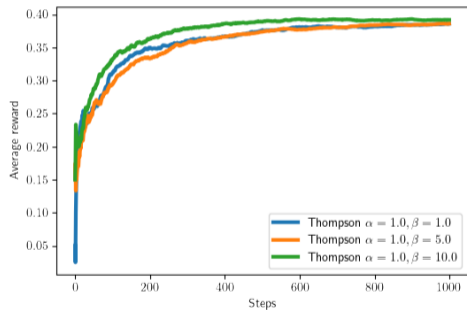
- **Action selection:**  $a^* \leftarrow \arg \max_a \theta_a$ .
- **Receive reward**  $R$  from arm  $a^*$ .
- **Update**  $S_a$  and  $F_a$ :
  - If  $R = 1$ ,  $S_a \leftarrow S_a + 1$ ;
  - If  $R = 0$ ,  $F_a \leftarrow F_a + 1$ .

---

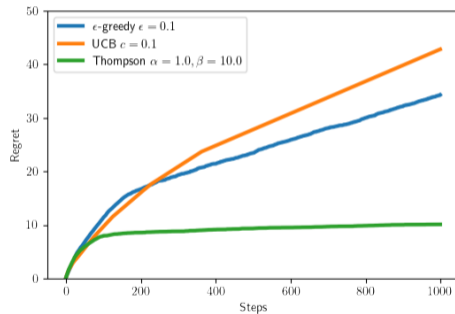
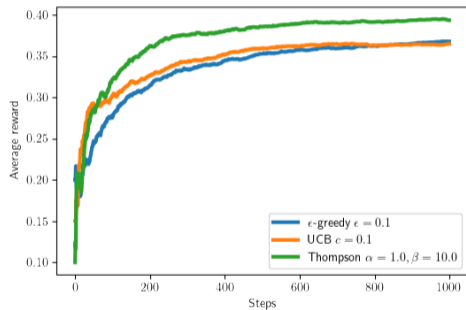
<sup>1</sup>The general algorithm can be found in Chapelle and Li [2011]



# Example: Four Bernoulli Arms (Revisited)



# Example: Four Bernoulli Arms (Cont'd)



# Theory

---

Define the **regret**  $T$ ,  $\mathcal{R}(T)$  as the difference between the cumulative reward under the **best action** and that under the **selected actions**, up to time  $T$ .

Theorem (UCB, Auer et al. [2002])

*The **expected regret** of the UCB algorithm  $\mathbb{E}\mathcal{R}(T)$  is upper bounded by  $C_1 \log(T)$  for some constant  $C_1 > 0$ .*

Theorem (TS, Agrawal and Goyal [2012])

*The **expected regret** of the Thompson sampling algorithm  $\mathbb{E}\mathcal{R}(T)$  is upper bounded by  $C_2 \log(T)$  for some constant  $C_2 > 0$ .*

- Both algorithms achieve logarithmic expected regret.
- Their performances are nearly the same as the oracle method that works as if the best action were known.
- $\epsilon$ -Greedy algorithm with a constant  $\epsilon$  has a **linear** expected regret (proportional to  $T$ ). More to discuss in seminar class.

1. Introduction and Course Overview

2. Multi-Armed Bandit

**3. Contextual Bandits**

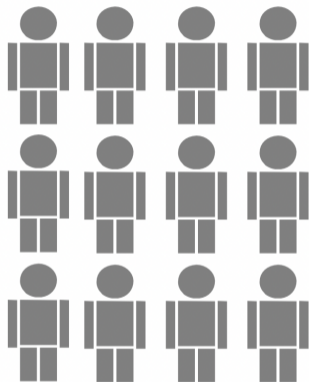
# Contextual Bandits

---

- Extension of MAB with **contextual** information.
- A **widely-used** model in medicine and technological industries.
- At time  $t$ , the agent
  - Observe a context  $S_t$ ;
  - Select an action  $A_t$ ;
  - Receives a reward  $R_t$  (depends on both  $S_t$  and  $A_t$ ).
- **Objective**: maximize cumulative reward.
- $\epsilon$ -greedy, **UCB** and **Thompson sampling** can be similarly adopted [see e.g., Chu et al., 2011, Agrawal and Goyal, 2013, Zhou et al., 2020, Zhang et al., 2020].

# Application I: Precision Medicine

---



**Patients**



**Treatment A**



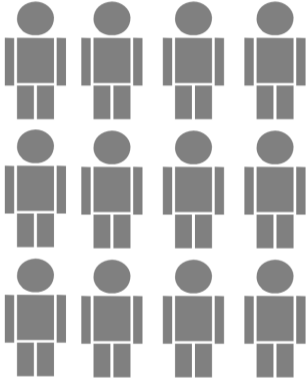
**Treatment B**



**Treatment C**

# One-Size-Fits-All

---

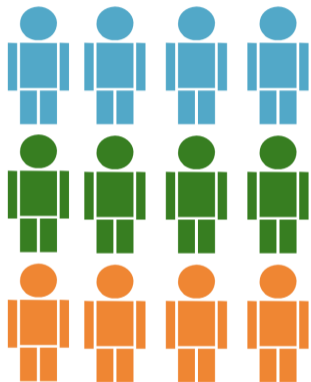


**Patients**



# Individualized Treatment Regime

---



**Patients**



**Treatment A**



**Treatment B**



**Treatment C**



# Application II: Personalized Recommendation



# Contextual Bandits Applications

---

- $S_t$ : Patient's or customer's baseline characteristics
- $A_t$ : Treatment (product) recommended to the patient (customer)
- $R_t$ : Patient's outcome or customer's action

# Summary

---

- Exploration-exploitation trade-off
- $\epsilon$ -greedy, UCB (the optimistic principle) and Thompson sampling
- Multi-armed bandits, contextual bandits

# Seminar Exercises

---

- Get started with **OpenAI Gym** ([link](#))
- Multi-armed bandits problem



# References I

---

- Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on learning theory*, pages 39–1. JMLR Workshop and Conference Proceedings, 2012.
- Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135. PMLR, 2013.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. *Advances in neural information processing systems*, 24:2249–2257, 2011.
- Elynn Y Chen, Rui Song, and Michael I Jordan. Reinforcement learning with heterogeneous data: Estimation and inference. *arXiv preprint arXiv:2202.00088*, 2022.

## References II

---

- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214. JMLR Workshop and Conference Proceedings, 2011.
- Luchen Li, Ignacio Albert-Smet, and Aldo A Faisal. Optimizing medical treatment for sepsis in intensive care: from reinforcement learning to pre-trial evaluation. *arXiv preprint arXiv:2003.06474*, 2020.
- Mengbing Li, Chengchun Shi, Zhenke Wu, and Piotr Fryzlewicz. Testing stationarity and change point detection in reinforcement learning. *arXiv preprint arXiv:2203.01707*, 2022.

## References III

---

- Peng Liao, Kristjan Greenewald, Predrag Klasnja, and Susan Murphy. Personalized heartsteps: A reinforcement learning algorithm for optimizing physical activity. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 4(1):1–22, 2020.
- Daniel J Lockett, Eric B Laber, Anna R Kahkoska, David M Maahs, Elizabeth Mayer-Davis, and Michael R Kosorok. Estimating dynamic treatment regimes in mobile health using v-learning. *Journal of the American Statistical Association*, 2019.
- Timothy NeCamp, Srijan Sen, Elena Frank, Maureen A Walton, Edward L Ionides, Yu Fang, Ambuj Tewari, and Zhenke Wu. Assessing real-time moderation for developing adaptive mobile health interventions for medical interns: Micro-randomized trial. *Journal of medical Internet research*, 22(3):e15033, 2020.

## References IV

---

- Chengchun Shi, Runzhe Wan, Rui Song, Wenbin Lu, and Ling Leng. Does the markov decision process fit the data: Testing for the markov property in sequential decision making. pages 8807–8817, 2020.
- Chengchun Shi, Sheng Zhang, Wenbin Lu, and Rui Song. Statistical inference of the value function for reinforcement learning in infinite-horizon settings. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 84(3):765–793, 2022.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. *arXiv preprint arXiv:2010.00827*, 2020.



## References V

---

- Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pages 11492–11502. PMLR, 2020.
- Wenzhuo Zhou, Ruoqing Zhu, and Annie Qu. Estimating optimal infinite horizon dynamic treatment regimes via pt-learning. *Journal of the American Statistical Association*, 0(0):1–14, 2022a.
- Yunzhe Zhou, Zhengling Qi, Chengchun Shi, and Lexin Li. Optimizing pessimism in dynamic treatment regimes: A bayesian learning approach. *arXiv preprint arXiv:2210.14420*, 2022b.

# Questions