

# **ST455: Reinforcement Learning**

## **Lecture 10: Offline Reinforcement Learning**

Chengchun Shi

# Lecture Outline

---

1. Introduction to Offline RL
2. The Pessimistic Principle
3. Model-based Offline Policy Optimization (MOPO)
4. An Overview of My Research

# Lecture Outline

---

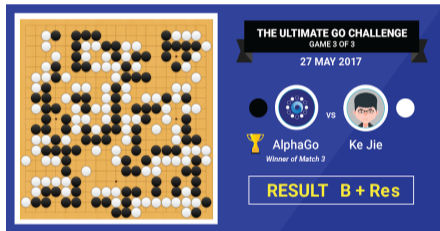
- 1. Introduction to Offline RL**
2. The Pessimistic Principle
3. Model-based Offline Policy Optimization (MOPO)
4. An Overview of My Research

# So Far, We Focused on Online RL Applications

---



(a) Video Games



(b) AlphaGo

# This Lecture Considers Offline Settings

---



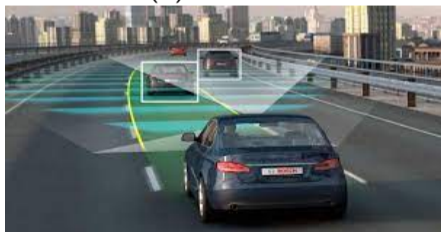
(a) Health Care



(b) Robotics



(c) Ridesharing



(d) Auto-driving

# This Lecture Considers Offline Settings (Cont'd)

---

- What is offline RL?
  - RL with a pre-collected historical dataset
- Why offline RL?
  - Online interaction with the environment is **impractical**
  - Either because online data collection is **expensive** (e.g., robotics or healthcare); rely on historical data
  - Or **dangerous** (e.g., healthcare, ridesharing or auto-driving)

# Online v.s. Offline RL

---

## Online RL:

- Data are **adaptively** generated, i.e., able to select **any** action at each time
- Data are **cheap** to generate, i.e., able to simulate **numerous** observations
- Likely to **satisfy** MDP assumption (Markovianity & time-homogeneity)

## Offline RL:

- Data are **pre-collected**, i.e., from an observational study
- Size of data is **limited**
- MDP assumption likely to be **violated** (Non-Markovianity or Non-stationarity)

# Offline RL Challenges and Solutions

---

- Data are **pre-collected**
  - Learning relies entirely on the historical data
  - Not possible to improve exploration
  - For actions that are less-explored, difficult to accurately learn their values
  - **Solution**: the pessimistic principle (focus of this lecture)
- Size of data is **limited**
  - **Solution**: develop sample-efficient RL algorithms (to be discussed in Lecture 11)
- **Violation** of MDP assumption
  - Cannot directly apply existing state-of-the-art RL algorithms
  - **Solution**: statistical hypothesis testing for model selection (to be covered in this lecture)



# Lecture Outline

---

1. Introduction to Offline RL
- 2. The Pessimistic Principle**
3. Model-based Offline Policy Optimization (MOPO)
4. An Overview of My Research

# Recap: Multi-Armed Bandit Problem

---

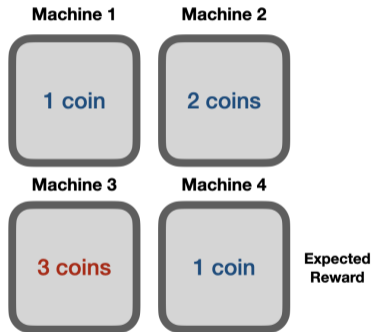


- The **simplest** RL problem
- A casino with **multiple** slot machines
- Playing each machine yields an independent **reward**.
- Limited knowledge (unknown reward distribution for each machine) and resources (**time**)
- **Objective**: determine which machine to pick at each time to maximize the expected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

# Offline Multi-Armed Bandit Problem

---

- $k$ -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**: Given  $\{A_t, R_t\}_{0 \leq t < T}$ , identify the best arm



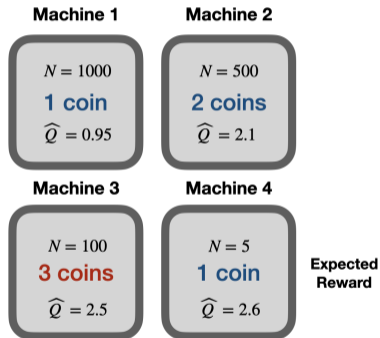
# Greedy Action Selection

- Action-value methods:

$$\hat{Q}(\mathbf{a}) = N^{-1}(\mathbf{a}) \sum_{t=0}^{T-1} R_t \mathbb{I}(\mathbf{A}_t = \mathbf{a})$$

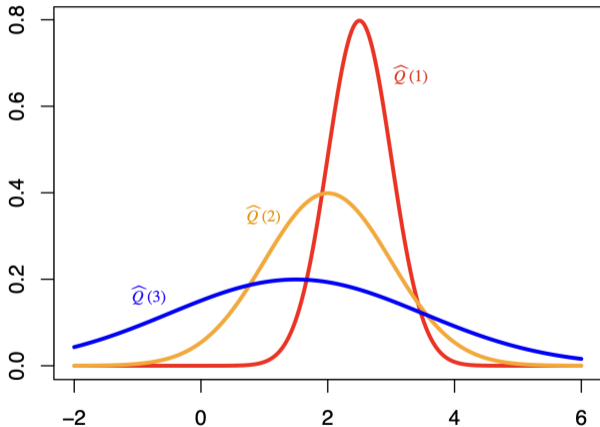
where  $N(\mathbf{a}) = \sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a})$   
denotes the action counter

- **Greedy policy:**  $\arg \max_{\mathbf{a}} \hat{Q}(\mathbf{a})$
- Less-explored action  $\rightarrow N(\mathbf{a})$  is small  
 $\rightarrow$  inaccurate  $\hat{Q}(\mathbf{a}) \rightarrow$  **suboptimal**  
policy (see the plot on the right)



# Recap: The Optimistic Principle

- Used in **online** settings to balance exploration-exploitation tradeoff
- 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
- Forms the basis for **upper confidence bound (UCB)**



# Recap: 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  $N_t(\mathbf{a})$  (number of times arm  $\mathbf{a}$  has been selected up to time  $t$ )
  - Large  $N_t(\mathbf{a}) \rightarrow$  small  $U_t(\mathbf{a})$ ;
  - Small  $N_t(\mathbf{a}) \rightarrow$  large  $U_t(\mathbf{a})$ .
- Select actions maximizing upper confidence bound

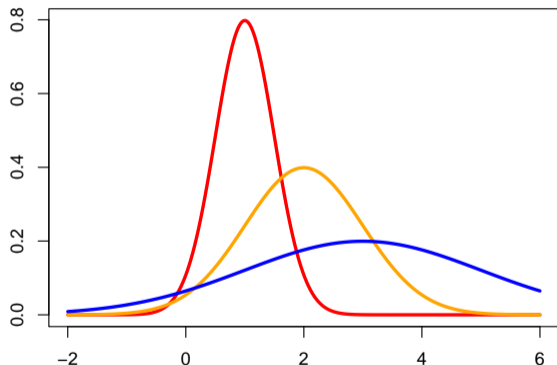
$$\mathbf{a}^* = \arg \max_{\mathbf{a}} [\hat{Q}_t(\mathbf{a}) + U_t(\mathbf{a})].$$

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

# The Pessimistic Principle

---

- In **offline** settings
- The less **uncertain** we are about an action-value
- The more **important** it is to use that action
- It could be the **best** action
- Likely to pick red action
- Yields the **lower confidence bound** (LCB) algorithm



# Lower Confidence Bound

---

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

$$Q(\mathbf{a}) \geq \hat{Q}(\mathbf{a}) - L(\mathbf{a}),$$

with high probability.

- $L(\mathbf{a})$  quantifies the **uncertainty** and depends on  $N(\mathbf{a})$  (number of times arm  $\mathbf{a}$  has been selected in the historical data)
  - Large  $N(\mathbf{a}) \rightarrow$  small  $L(\mathbf{a})$ ;
  - Small  $N(\mathbf{a}) \rightarrow$  large  $L(\mathbf{a})$ .
- Select actions maximizing lower confidence bound

$$\mathbf{a}^* = \arg \max_{\mathbf{a}} [\hat{Q}(\mathbf{a}) - L(\mathbf{a})].$$



## Lower Confidence Bound (Cont'd)

---

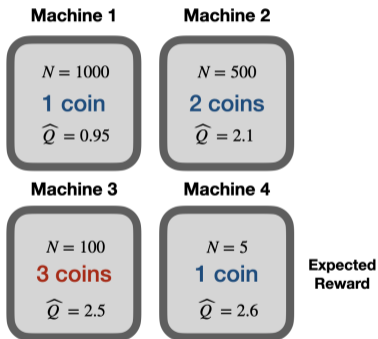
- Set  $L(\mathbf{a}) = \sqrt{c \log(\mathcal{T}) / N(\mathbf{a})}$  for some positive constant  $c$  where  $\mathcal{T}$  is the sample size of historical data
- According to **Hoeffding's inequality** ([link](#)), when rewards are bounded between  $0$  and  $1$ , the event

$$|Q(\mathbf{a}) - \hat{Q}(\mathbf{a})| \leq L(\mathbf{a}),$$

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

# Lower Confidence Bound (Cont'd)

- $\hat{Q}(4) > \hat{Q}(3)$
- $T = 1605$ . Set  $c = 1$
- $L(3) = \sqrt{\log(T)/N(3)} = 0.272$
- $L(4) = \sqrt{\log(T)/N(4)} = 1.215$
- $\hat{Q}(3) - L(3) > \hat{Q}(4) - L(4)$
- $\hat{Q}(3) - L(3) > \max(\hat{Q}(1), \hat{Q}(2))$
- Correctly identify optimal action



# Algorithm

---

- **Input:** some positive constant  $c$ , offline data  $\{(A_t, R_t)\}_{0 \leq t < T}$ .
- **Initialization:**  $t = 0$ ,  $\hat{Q}(a) = 0$ ,  $N(a) = 0$ , for  $a = 1, 2, \dots, k$ .
- **While**  $t < T$ :
  - Update  $N$ :  $N(A_t) \leftarrow N(A_t) + 1$ .
  - Update  $\hat{Q}$ :

$$\hat{Q}(A_t) \leftarrow \frac{N(A_t) - 1}{N(A_t)} \hat{Q}(A_t) + \frac{1}{N(A_t)} R_t.$$

- Update  $t$ :  $t \leftarrow t + 1$ .
- **LCB action selection:**

$$a^* \leftarrow \arg \max_a [\hat{Q}(a) - \sqrt{c \log(T) / N(a)}].$$

# Theory

---

Define the regret, as the difference between the expected reward under the **best arm** and that under the **selected arm**.

## Theorem (Greedy Action Selection)

*Regret of greedy action selection is upper bounded by  $2 \max_{\mathbf{a}} |\hat{Q}(\mathbf{a}) - Q(\mathbf{a})|$ , whose value is bounded by  $2\sqrt{c \log(\mathbf{T}) / \min_{\mathbf{a}} \mathbf{N}(\mathbf{a})}$  (according to Hoeffding's inequality) with probability approaching 1*

- The upper bound depends on the estimation error of **each** Q-estimator
- The regret is small when **each** arm has sufficiently many observations
- However, it would yield a large regret when one arm is **less-explored**
- This reveals the **limitation** of greedy action selection
- Proof is simple (see Appendix)

# Theory (Cont'd)

---

Theorem (LCB; see also Jin et al. [2021])

*Regret of the LCB algorithm is upper bounded by  $2\sqrt{c \log(\mathbf{T})/\mathbf{N}(\mathbf{a}^{opt})}$  where  $\mathbf{a}^{opt}$  denotes the best arm with probability approaching **1***

- The upper bound depends on the estimation error of best arm's Q-estimator **only**
- The regret is small when the **best** arm has sufficiently many observations
- This is much weaker than requiring **each** arm to have sufficiently many observations
- This reveals the **advantage** of LCB algorithm
- Proof given in the Appendix

# Lecture Outline

---

1. Introduction to Offline RL
2. The Pessimistic Principle
- 3. Model-based Offline Policy Optimization (MOPO)**
4. An Overview of My Research

# Offline RL and Fitted Q-Iteration

---

- Offline data:  $\{(S_t, A_t, R_t) : 0 \leq t \leq T\}$
- Fitted Q-Iteration can be naturally applied by repeating
  1. Compute  $\hat{Q}$  as the argmin of

$$\arg \min_Q \sum_t \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \right]^2$$

2. Set  $\tilde{Q} = \hat{Q}$
- **Limitation:** for less-explored state-action pairs, their Q-values **cannot** be learned accurately
  - **Solution:** the pessimistic principle

# Pessimistic Principle in RL

---

- In multi-armed bandit, we select action to maximize lower confidence bound

$$\mathbf{a}^* = \arg \max_{\mathbf{a}} [\widehat{Q}(\mathbf{a}) - L(\mathbf{a})]$$

- In more general RL, we can adopt a similar principle by setting

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} \mathbf{1}, & \text{if } \mathbf{a} = \arg \max \widehat{Q}(\mathbf{a}, \mathbf{s}) - L(\mathbf{a}, \mathbf{s}) \\ \mathbf{0}, & \text{otherwise} \end{cases}$$

where the lower bound satisfies that with probability approaching  $\mathbf{1}$ ,

$$Q^{\pi^{\text{opt}}}(\mathbf{a}, \mathbf{s}) \geq \widehat{Q}(\mathbf{a}, \mathbf{s}) - L(\mathbf{a}, \mathbf{s}), \quad \forall \mathbf{a}, \mathbf{s}.$$

- Many offline algorithms [see e.g., Wu et al., 2019, Kumar et al., 2020, Levine et al., 2020] adopt similar ideas, but do not directly use the above formula



# Model-based Offline Policy Optimisation (MOPO)

---

- As we discussed in Lecture 9, **model-based** method is preferred in offline settings
- Online RL algorithms are **not** applicable, as adaptive interaction is not feasible
- Model-based method
  - learns a model using the **offline** data
  - allows to **adaptively** generate data based on the model
  - applies **online** RL algorithms to simulated data for policy optimisation
  - embraces the power of online RL algorithms for offline policy optimisation
- MOPO [Yu et al., 2020] integrates model-based method with **pessimistic** principle

# MOPO: Offline Model Learning

---

- Learn the conditional distribution of  $(\mathbf{S}_{t+1}, \mathbf{R}_t)$  given  $(\mathbf{A}_t, \mathbf{S}_t)$
- Approximate the conditional distribution using Gaussian, i.e.,

$$(\mathbf{S}_{t+1}, \mathbf{R}_t) | (\mathbf{A}_t, \mathbf{S}_t) \sim N(\boldsymbol{\mu}_\theta(\mathbf{A}_t, \mathbf{S}_t), \boldsymbol{\Sigma}_\phi(\mathbf{A}_t, \mathbf{S}_t))$$

- Parametrize  $\boldsymbol{\mu}_\theta$  and  $\boldsymbol{\Sigma}_\phi$  using e.g., neural networks
- Use bootstrap to learn  $N$  different models  $\{\mathcal{M}_i\}_{i=1, \dots, N}$

# MOPO: The Pessimism Principle

---

- Penalize reward to incorporate pessimism
- Simulate reward  $r$  given the state-action pair  $(s, a)$  from model
- Define the **transformed** reward

$$\tilde{r} = r - L(a, s),$$

for some lower bound  $L(a, s)$  that quantifies the **uncertainty** of model

- More uncertain  $\rightarrow$  smaller transformed reward
- Less uncertain  $\rightarrow$  larger transformed reward
- Apply online RL to transformed data (see next slide)

# MOPO: Adaptive Data Simulation

---

## 1. Action simulation

- For **value-based** method, sample actions using  $\epsilon$ -greedy policy
- For **policy-gradient** method, sample actions using the estimated policy

## 2. Reward and next-state simulation

- Randomly pick a model  $\mathcal{M}_i = N(\mu_{\theta_i}(\mathbf{A}_t, \mathbf{S}_t), \Sigma_{\phi_i}(\mathbf{A}_t, \mathbf{S}_t))$
- Sample  $(\mathbf{S}_{t+1}, \mathbf{R}_t)$  from this Gaussian model
- Compute transformed reward  $\widetilde{\mathbf{R}}_t = \mathbf{R}_t - L(\mathbf{A}_t, \mathbf{S}_t)$
- Use  $(\mathbf{S}_t, \mathbf{A}_t, \widetilde{\mathbf{R}}_t, \mathbf{S}_{t+1})$  to update the policy/Q-function

## 3. Repeat the above two steps for data simulation and policy learning

# MOPO: Pseudocode

---

---

**Algorithm 2** MOPO instantiation with regularized probabilistic dynamics and ensemble uncertainty

---

**Require:** reward penalty coefficient  $\lambda$  rollout horizon  $h$ , rollout batch size  $b$ .

- 1: Train on batch data  $\mathcal{D}_{\text{env}}$  an ensemble of  $N$  probabilistic dynamics  $\{\widehat{T}^i(s', r | s, a) = \mathcal{N}(\mu^i(s, a), \Sigma^i(s, a))\}_{i=1}^N$ .
  - 2: Initialize policy  $\pi$  and empty replay buffer  $\mathcal{D}_{\text{model}} \leftarrow \emptyset$ .
  - 3: **for** epoch  $1, 2, \dots$  **do** ▷ This for-loop is essentially one outer iteration of MBPO
  - 4:     **for**  $1, 2, \dots, b$  (in parallel) **do**
  - 5:         Sample state  $s_1$  from  $\mathcal{D}_{\text{env}}$  for the initialization of the rollout.
  - 6:         **for**  $j = 1, 2, \dots, h$  **do**
  - 7:             Sample an action  $a_j \sim \pi(s_j)$ .
  - 8:             Randomly pick dynamics  $\widehat{T}$  from  $\{\widehat{T}^i\}_{i=1}^N$  and sample  $s_{j+1}, r_j \sim \widehat{T}(s_j, a_j)$ .
  - 9:             Compute  $\tilde{r}_j = r_j - \lambda \max_{i=1}^N \|\Sigma^i(s_j, a_j)\|_F$ .
  - 10:             Add sample  $(s_j, a_j, \tilde{r}_j, s_{j+1})$  to  $\mathcal{D}_{\text{model}}$ .
  - 11:     Drawing samples from  $\mathcal{D}_{\text{env}} \cup \mathcal{D}_{\text{model}}$ , use SAC to update  $\pi$ .
-

# Lecture Outline

---

1. Introduction to Offline RL
2. The Pessimistic Principle
3. Model-based Offline Policy Optimization (MOPO)
- 4. An Overview of My Research**

# Recap: The Agent's Policy

---



# Recap: Foundations of RL

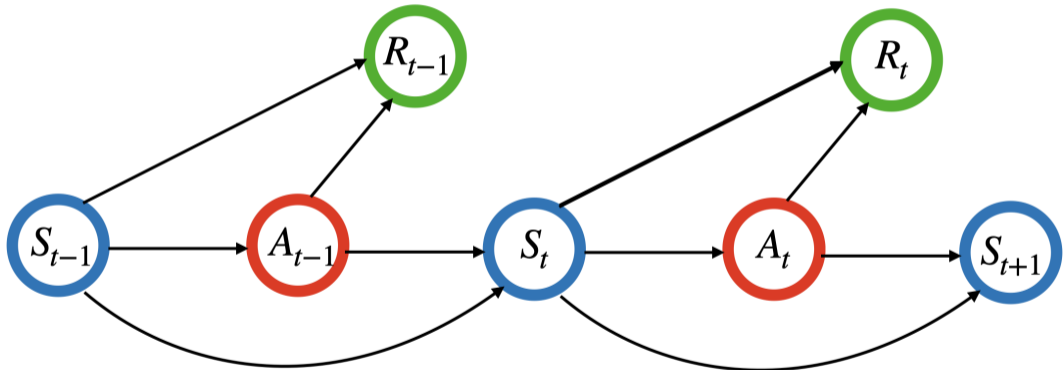
---

- **RL algorithms:** policy iteration, value iteration (Lecture 3), SARSA, Q-learning (Lecture 4), gradient-based methods, fitted Q-iteration (Lecture 5), deep Q-network (Lecture 7), REINFORCE, actor critic (Lecture 8), Dyna-Q (Lecture 9)
- **Foundations** of aforementioned algorithms:
  - **Markov decision process** [MDP, Puterman, 2014]: ensures the optimal policy is **stationary**, and is **not** history-dependent
  - **Markov assumption:** conditional on the present (e.g.,  $S_t, A_t$ ), the future (e.g.,  $R_t, S_{t+1}$ ) and the past data history are independent
  - **Time-homogeneity assumption:** The conditional distribution of  $(R_t, S_{t+1})$  given  $(S_t = s, A_t = a)$  is time-homogeneous



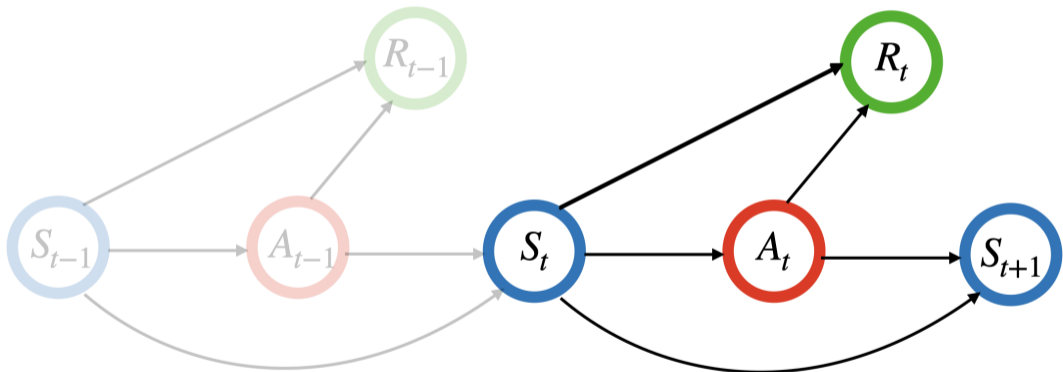
## Recap: Markov Assumption

---



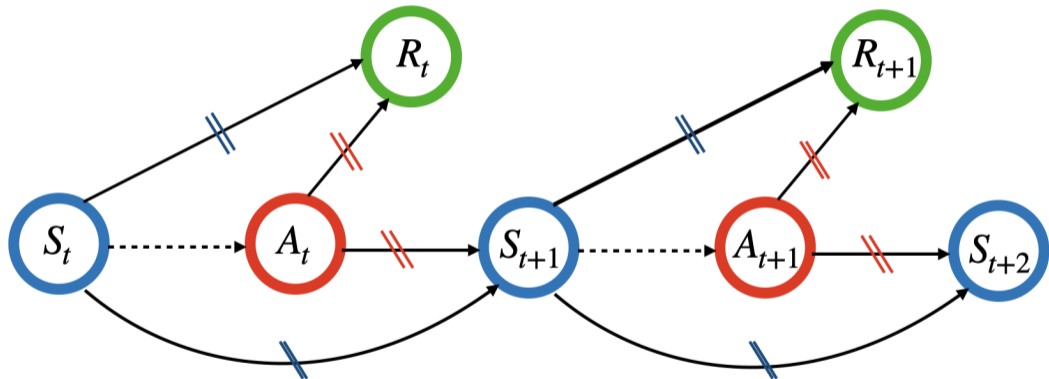
## Recap: Markov Assumption

---



# Recap: Time-Homogeneity Assumption

---



# Violation of MDP Assumption

---

- Violation of Markov assumption
  - Statistical hypothesis testing for model selection: MDP, high-order MDP ( $k$ th order for  $k \geq 2$ ), POMDP ( $\infty$ th order MDP)
- Violation of time-homogeneity assumption
  - Statistical hypothesis testing for selecting the “best data segment”

# Markov and Non-Markov Models

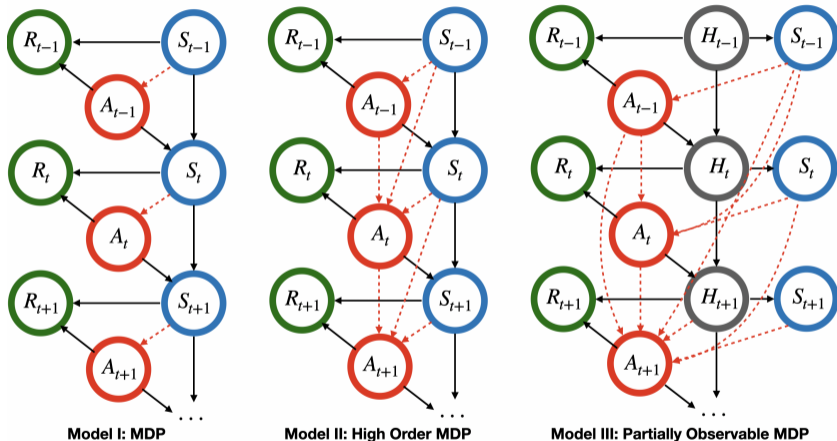


Figure: Causal diagrams for MDPs, HMDPs and POMDPs. The solid lines represent the causal relationships and the dashed lines indicate the information needed to implement the optimal policy.  $\{H_t\}_t$  denotes latent variables.

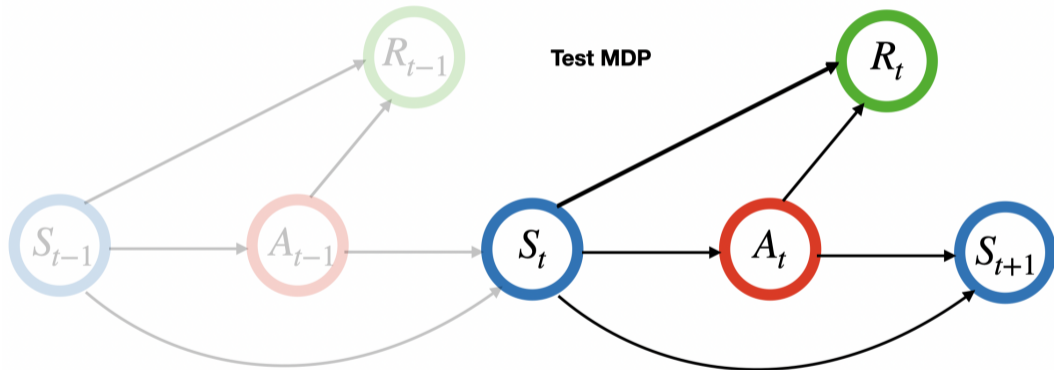
# Test Markov Assumption [Shi et al., 2020]

---

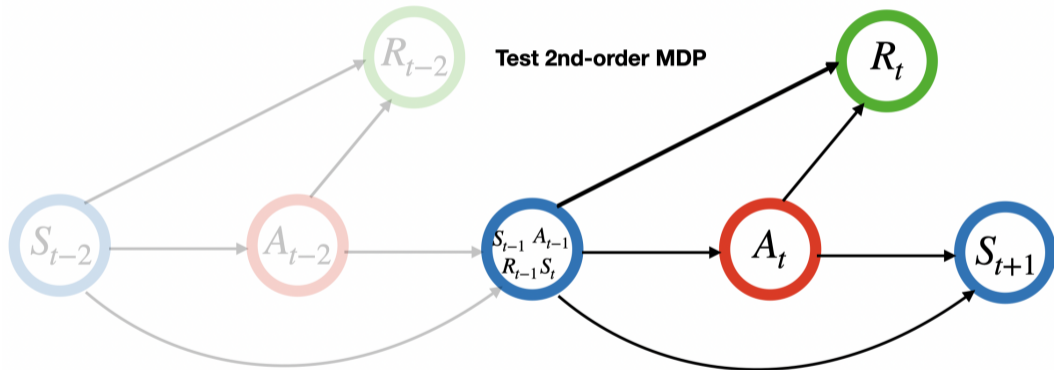
- Develop a **forward-backward learning procedure** to test the Markov assumption (MA) in RL
  - Null hypothesis  $\mathcal{H}_0$ : MA holds (MDP)
  - Alternative hypothesis  $\mathcal{H}_1$ : MA is violated (high-order MDP, POMDP)
- Sequentially apply the test for **model selection**
  - Suppose the data follows a  $K$ th order MDP
  - Sequentially test whether it is  $k$ th order for  $k = 1, 2, \dots$
  - by concatenating  $S_t$  with  $\{(S_{t-j}, A_{t-j}, R_{t-j})\}$  for  $1 \leq j < k$
  - $\mathcal{H}_0$  holds when  $k \geq K$  and  $\mathcal{H}_1$  holds otherwise
  - Select the model when  $\mathcal{H}_0$  is not rejected for the first time

# Test Markov Assumption (Cont'd)

---

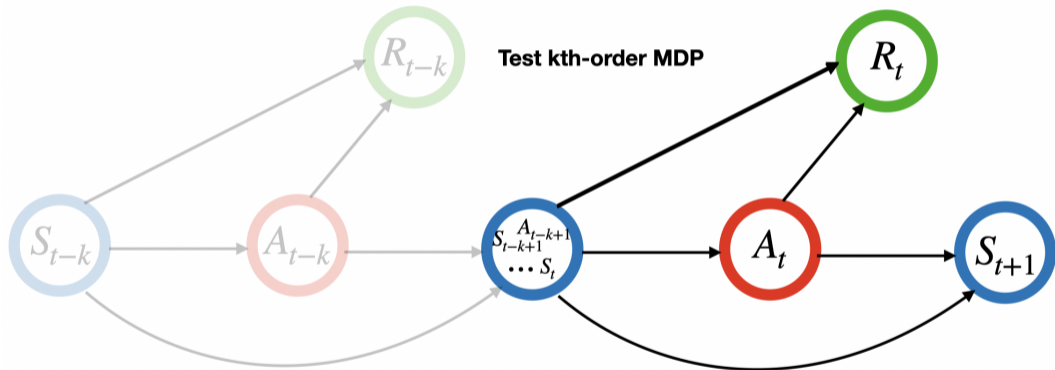


# Test Markov Assumption (Cont'd)





# Test Markov Assumption (Cont'd)



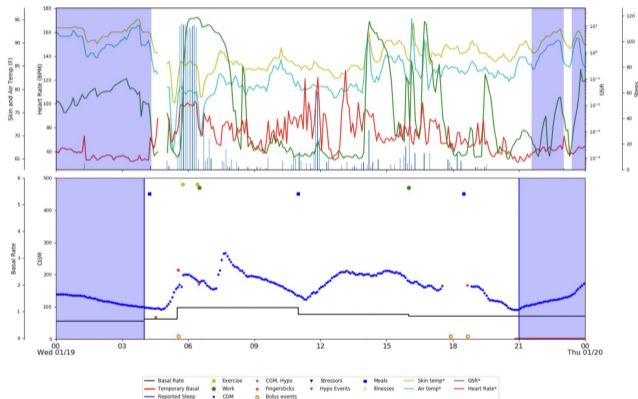
# Test Markov Assumption (Cont'd)

---

- Uncritical to **online** domains:
  - Try different models online and see which model yields the best reward
- Critical to **offline** domains:
  - $K$  remains unknown without prior knowledge
  - Cannot adaptively generate data
  - For **under-fitted** models ( $k < K$ ), any stationary policy is not optimal
  - For **over-fitted** models ( $k > K$ ), the estimated policy might be very noisy due to the inclusion of many irrelevant lagged variables

# Diabetes

- Management of **Type-I diabetes**
- **Subject:** Patients with diabetes.
- **Objective:** Develop treatment policy to determine whether patients need to inject insulin at each time to improve their health
- $S_t$ : Patient's **glucose levels, food intake, exercise intensity**
- $A_t$ : **Insulin doses** injected
- $R_t$ : **Index of Glycemic Control** (function of patient's glucose level)



# Diabetes (Cont'd)

---

- **Analysis I:**

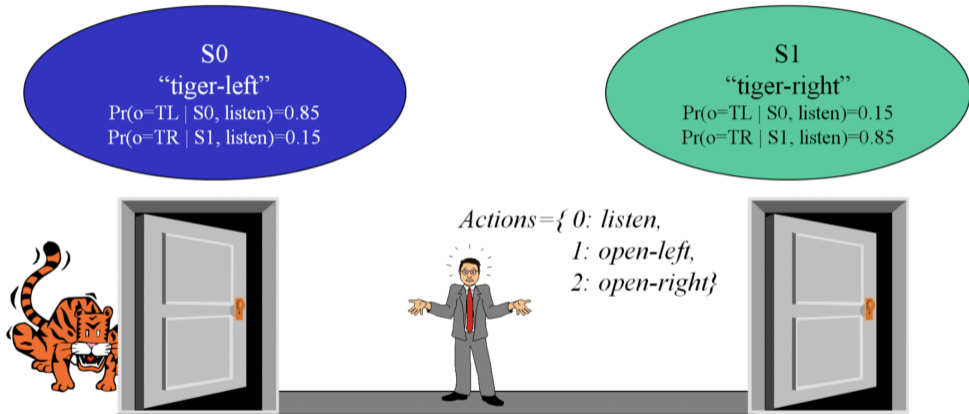
- sequentially apply our test to determine the order of MDP
- conclude it is a **fourth-order** MDP

- **Analysis II:**

- split the data into training/testing samples
- policy optimization based on **fitted-Q iteration**, by assuming it is a  $k$ -th order MDP for  $k = 1, \dots, 10$
- policy evaluation based on **fitted-Q evaluation** (to be covered in Lecture 11)
- use **random forest** to model the Q-function
- repeat the above procedure to compute the average value of policies computed under each MDP model assumption

order	1	2	3	4	5	6	7	8	9	10
value	-90.8	-57.5	-63.8	<b>-52.6</b>	-56.2	-60.1	-63.7	-54.9	-65.1	-59.6

# Tiger



## Reward Function

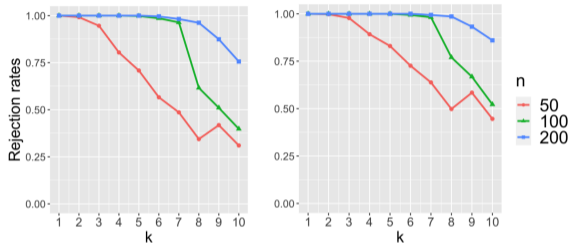
- Penalty for wrong opening: -100
- Reward for correct opening: +10
- Cost for listening action: -1

## Observations

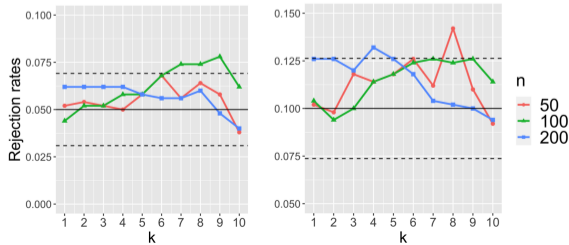
- to hear the tiger on the left (TL)
- to hear the tiger on the right (TR)

# Tiger (Cont'd)

- Under the alternative hypothesis (MA is violated).  $\alpha = (0.05, 0.1)$  from left to right.



- Under the null hypothesis (MA holds).  $\alpha = (0.05, 0.1)$  from left to right.



# Test Time-Homogeneity [Li et al., 2022]

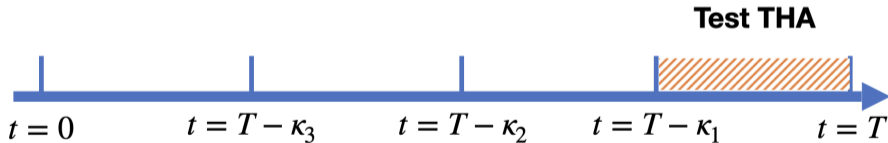
---

- Under **time-inhomogeneity**, using all data is not reasonable
- Natural to use **more recent observations** for policy optimisation
- **Challenging** to select the **best data “segment”**
  - Including too many past observations yields a suboptimal policy
  - Using only a few recent observations results in a very noisy policy
- Develop a **test procedure** for the time-homogeneity assumption (THA) in RL
  - Null hypothesis  $\mathcal{H}_0$ : THA holds (MDP)
  - Alternative hypothesis  $\mathcal{H}_1$ : THA is violated (Time-Varying MDP)
- Sequentially apply the test for selecting the **best data “segment”**

# Test Time-Homogeneity (Cont'd)

---

- Sequentially apply the test for selecting the **best data “segment”**
  - Sequentially test whether THA holds on the data interval  $[T - \kappa, T]$  for  $\kappa_1 < \kappa_2 < \kappa_3 < \dots$
  - Suppose THA is first rejected at some  $\kappa = \kappa_{j_0}$
  - Use the data subset within the interval  $[T - \kappa_{j_0-1}, T]$  for policy optimisation



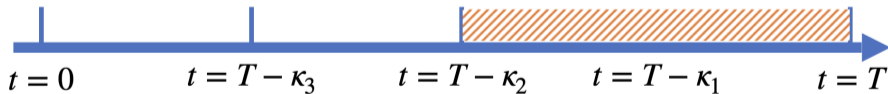


# Test Time-Homogeneity (Cont'd)

---

- Sequentially apply the test for selecting the **best data “segment”**
  - Sequentially test whether THA holds on the data interval  $[T - \kappa, T]$  for  $\kappa_1 < \kappa_2 < \kappa_3 < \dots$
  - Suppose THA is first rejected at some  $\kappa = \kappa_{j_0}$
  - Use the data subset within the interval  $[T - \kappa_{j_0-1}, T]$  for policy optimisation

**Not rejected. Combine more data**

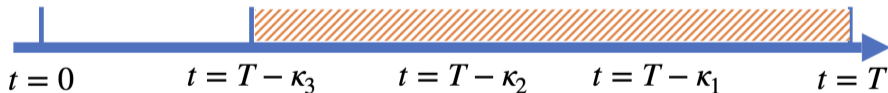


# Test Time-Homogeneity (Cont'd)

---

- Sequentially apply the test for selecting the **best data “segment”**
  - Sequentially test whether THA holds on the data interval  $[T - \kappa, T]$  for  $\kappa_1 < \kappa_2 < \kappa_3 < \dots$
  - Suppose THA is first rejected at some  $\kappa = \kappa_{j_0}$
  - Use the data subset within the interval  $[T - \kappa_{j_0-1}, T]$  for policy optimisation

**Not rejected. Combine more data**

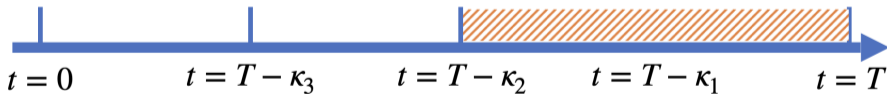


# Test Time-Homogeneity (Cont'd)

---

- Sequentially apply the test for selecting the **best data “segment”**
  - Sequentially test whether THA holds on the data interval  $[T - \kappa, T]$  for  $\kappa_1 < \kappa_2 < \kappa_3 < \dots$
  - Suppose THA is first rejected at some  $\kappa = \kappa_{j_0}$
  - Use the data subset within the interval  $[T - \kappa_{j_0-1}, T]$  for policy optimisation

**Rejected. Use the last data interval**



# Intern Health Study

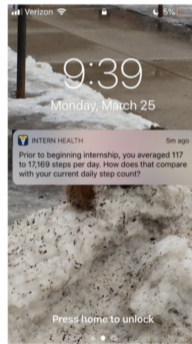
- **Subject:** First-year medical interns
- **Objective:** Develop treatment policy to determine whether to send certain text messages to interns to improve their health
- $S_t$ : Interns' mood scores, sleep hours and step counts
- $A_t$ : Send text notifications or not
- $R_t$ : Step counts



(i) App Dashboard

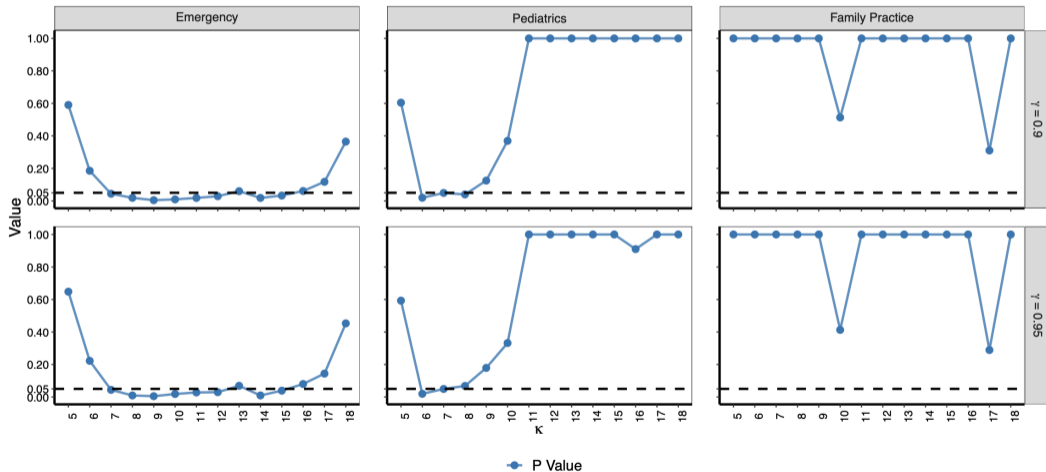


(ii) Mood EMA



(iii) Notifications

# Intern Health Study (Cont'd)



# Intern Health Study (Cont'd)

Number of Change Points	Specialty	Method	$\gamma = 0.9$	$\gamma = 0.95$
$\geq 1$	Emergency	Proposed	8237.16	8295.99
		Overall	8108.13	8127.55
		Behavior	7823.75	7777.32
		Random	8114.78	8080.27
$\geq 2$	Pediatrics	Proposed	7883.08	7848.57
		Overall	7925.44	7960.12
		Behavior	7730.98	7721.29
		Random	7807.52	7815.30
0	Family Practice	Proposed	8062.50	7983.69
		Overall	8062.50	7983.69
		Behavior	7967.67	7957.24
		Random	7983.52	7969.31

TABLE 3

*Mean value estimates using decision tree in analysis of IHS. Values are normalised by multiplying  $1 - \gamma$ . All values are evaluated over 10 splits of data.*

- Mean value is the weekly average step counts per day
- The proposed method improves mean value by 50 – 150 steps, compared to the behavior policy

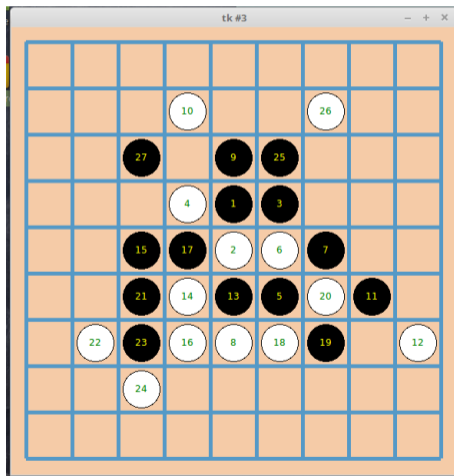
# Summary

---

- Offline RL v.s. online RL
- The pessimistic principle
- Lower confidence bound
- Model-based offline policy optimisation
- Statistical hypothesis testing

# Seminar Exercise

- Solutions to HW9 (Deadline: Wed 12 pm)
- Implementation of AlphaZero on Gomoku





# References I

---

- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Mengbing Li, Chengchun Shi, Zhenke Wu, and Piotr Fryzlewicz. Reinforcement learning in possibly nonstationary environments. *arXiv preprint arXiv:2203.01707*, 2022.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

## References II

---

- 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. *arXiv preprint arXiv:2002.01751*, 2020.
- Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing Systems*, 33:14129–14142, 2020.

# Questions

# Appendix: Proof of Regret

---

Consider the regret of greedy action selection first. Let  $\mathbf{a}^*$  denote the action selected by the greedy policy. By definition, the regret is given by  $Q(\mathbf{a}^{opt}) - Q(\mathbf{a}^*)$ . Notice that

$$\begin{aligned} Q(\mathbf{a}^{opt}) - Q(\mathbf{a}^*) &= Q(\mathbf{a}^{opt}) - \widehat{Q}(\mathbf{a}^{opt}) + \widehat{Q}(\mathbf{a}^{opt}) - \widehat{Q}(\mathbf{a}^*) + \widehat{Q}(\mathbf{a}^*) - Q(\mathbf{a}^*) \\ &\leq Q(\mathbf{a}^{opt}) - \widehat{Q}(\mathbf{a}^{opt}) + \widehat{Q}(\mathbf{a}^*) - Q(\mathbf{a}^*), \end{aligned}$$

as  $\mathbf{a}^*$  maximizes  $\arg \max_{\mathbf{a}} \widehat{Q}(\mathbf{a})$  by definition.

It is immediate to see that the right-hand-side is upper bounded by  $2 \max_{\mathbf{a}} |\widehat{Q}(\mathbf{a}) - Q(\mathbf{a})|$ . The proof is thus completed.

## Appendix: Proof of Regret (Cont'd)

---

Next, consider the regret of the LCB algorithm. Let  $\mathbf{a}^*$  denote the action selected by the LCB algorithm. By definition of  $L(\mathbf{a}^*)$ , we have with probability approaching  $\mathbf{1}$  that

$$Q(\mathbf{a}^{opt}) - Q(\mathbf{a}^*) \leq Q(\mathbf{a}^{opt}) - \widehat{Q}(\mathbf{a}^*) + L(\mathbf{a}^*).$$

According to the LCB algorithm,  $\widehat{Q}(\mathbf{a}^*) - L(\mathbf{a}^*) \geq \widehat{Q}(\mathbf{a}^{opt}) - L(\mathbf{a}^{opt})$ . It follows that the right-hand-side is upper bounded by

$$Q(\mathbf{a}^{opt}) - \widehat{Q}(\mathbf{a}^{opt}) + L(\mathbf{a}^{opt}),$$

which is further bounded by  $2L(\mathbf{a}^{opt})$ , by definition. The proof is completed by directly applying Hoeffding's inequality.