

# Reinforcement Learning

Lecture 3: Q-Learning and Beyond

Chengchun Shi

# Lecture Outline

---

1. Tabular Q-Learning
2. Fitted Q-Iteration
3. Case Study I: Deep Q-Network (DQN) in Atari
4. Case Study II: TD Learning in Ridesharing Platforms

# Lecture Outline

---

1. Tabular Q-Learning
2. Fitted Q-Iteration
3. Case Study I: Deep Q-Network (DQN) in Atari
4. Case Study II: TD Learning in Ridesharing Platforms

# Introduction to Q-Learning

---

- One of the most **popular** class of RL algorithms
- Variants include double Q-learning, fitted Q-iteration, deep Q-network (DQN), quantile DQN
- **Main idea:** learn the optimal Q-function  $Q^{\pi^{\text{opt}}}$  based on the Bellman **optimality** equation and derive the optimal policy

$$\pi^{\text{opt}}(\textcolor{blue}{s}) = \arg \max_{\textcolor{red}{a}} Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})$$

- Focus on **tabular Q-learning** (finite MDP, discrete state and action)

$$\pi^{\text{opt}}(\textcolor{blue}{s}) = \arg \max_{\textcolor{red}{a}} Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})?$$

---

- $Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})$  is the **value** of the policy that
  - Assigns  $\textcolor{red}{a}$  at the initial decision time;
  - Follow  $\pi^{\text{opt}}$  afterwards
- $Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \pi^{\text{opt}}(\textcolor{blue}{s})) = V^{\pi^{\text{opt}}}(\textcolor{blue}{s})$  is the value under the optimal policy  $\pi^{\text{opt}}$
- $\pi^{\text{opt}}$  is stationary and is no worse than any **history-dependent** policy (Lecture 1)

$$Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a}) \leq V^{\pi^{\text{opt}}}(\textcolor{blue}{s}) = Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \pi^{\text{opt}}(\textcolor{blue}{s})), \quad \forall \textcolor{red}{a}.$$

- It follows that

$$\pi^{\text{opt}}(\textcolor{blue}{s}) = \arg \max_{\textcolor{red}{a}} Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})$$

# Bellman Optimality Equation

---

- Bellman optimality equation for the optimal value function:

$$V^{\pi^{\text{opt}}}(s) = \max_a \mathbb{E}[R_t + \gamma V^{\pi^{\text{opt}}}(S_{t+1}) | A_t = a, S_t = s].$$

- Bellman optimality equation for the optimal Q-function:

$$Q^{\pi^{\text{opt}}}(s, a) = \mathbb{E} \left[ R_t + \gamma \max_{a'} Q^{\pi^{\text{opt}}}(S_{t+1}, a') | A_t = a, S_t = s \right].$$

- **Proof:** according to Bellman equation,

$$Q^{\pi^{\text{opt}}}(s, a) = \mathbb{E} \left[ R_t + \gamma Q^{\pi^{\text{opt}}}(S_{t+1}, \pi^{\text{opt}}(S_{t+1})) | A_t = a, S_t = s \right]$$

This together with  $\pi^{\text{opt}}(s) = \arg \max_a Q^{\pi^{\text{opt}}}(s, a)$  yields the result.

# Q-Learning: an Off-Policy TD Control

---

- One-step SARSA update:

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha [R_t + \gamma Q(\mathbf{S}_{t+1}, \mathbf{A}_{t+1}) - Q(\mathbf{S}_t, \mathbf{A}_t)]$$

- One-step tabular Q-learning update:

$$Q(\mathbf{S}_t, \mathbf{A}_t) \leftarrow Q(\mathbf{S}_t, \mathbf{A}_t) + \alpha \left[ R_t + \gamma \max_{\mathbf{a}} Q(\mathbf{S}_{t+1}, \mathbf{a}) - Q(\mathbf{S}_t, \mathbf{A}_t) \right]$$

- In Q-learning, the action in the target is independent of the **behavior policy**
- The behavior policy has an effect on which state-actions are visited

# Tabular Q-Learning: Pseudocode

---

- **Initialization:**  $Q$  arbitrary
- **Repeat** for each episode:

**Initialize** state  $s$

**Repeat** for each step of the episode:

$a \leftarrow$  action from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

        Take action  $a$ , observe reward  $r$  and next state  $s'$

$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$

$s \leftarrow s'$

**until**  $s$  is a terminal state

# On-Policy v.s Off-Policy

---

- Q-learning is **off-policy**:
  - Updates Q-values using Q-value of next state  $s'$  and **greedy** action  $a'$
  - Q-value evaluates the value of **greedy policy** instead of **current policy**
- SARSA is **on-policy**:
  - Updates Q-values using Q-value of next state  $s'$  and **current policy's** action  $a'$
  - Q-value evaluates the value of **current policy**

# Recap: Convergence of SARSA

## Theorem

SARSA converges to the optimal  $Q$ -function,  $\mathbf{Q}(\mathbf{s}, \mathbf{a}) \rightarrow \mathbf{Q}^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})$  for any  $\mathbf{s}$  and  $\mathbf{a}$ , if

- All state-action pairs are explored infinitely many times,

$$\sum_{t=0}^{\infty} \mathbb{I}(\mathbf{S}_t = \mathbf{s}, \mathbf{A}_t = \mathbf{a}) = \infty.$$

- The policy converges to a greedy policy,

$$\lim_{t \rightarrow \infty} \pi_t(\mathbf{a}|\mathbf{s}) = \mathbb{I}(\mathbf{a} = \arg \max_{\mathbf{a}'} \mathbf{Q}_t(\mathbf{s}, \mathbf{a}'))$$

- Robbins-Monro sequence of step-sizes [Robbins and Monro, 1951],

$$\sum_{t=0}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

# Convergence of Tabular Q-Learning

Theorem (Melo [2001])

*Q-learning converges to the optimal Q-function if*

- All state-action pairs are explored infinitely many times,

$$\sum_{t=0}^{\infty} \mathbb{I}(S_t = s, A_t = a) = \infty.$$

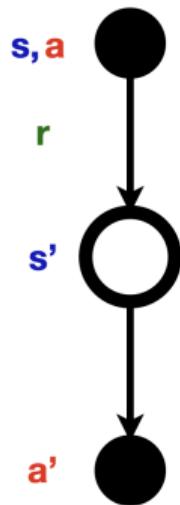
- Robbins-Monro sequence of step-sizes [Robbins and Monro, 1951],

$$\sum_{t=0}^{\infty} \alpha_t = \infty \text{ and } \sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

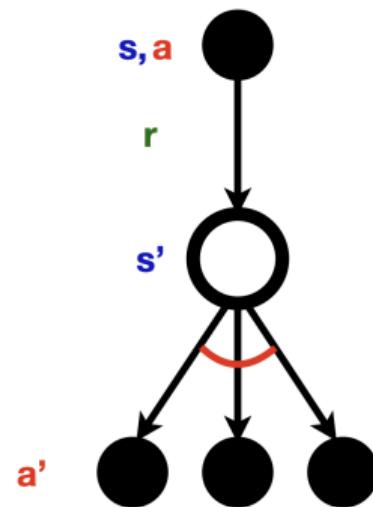
- Only requires  $\epsilon$  to be strictly positive
- No need to require  $\epsilon$  to decay to zero
- Q-learning converges even if the behavior policy is far from the optimal

# Backup Diagram

---



SARSA



Q-Learning

# Cliff Walking Example

---

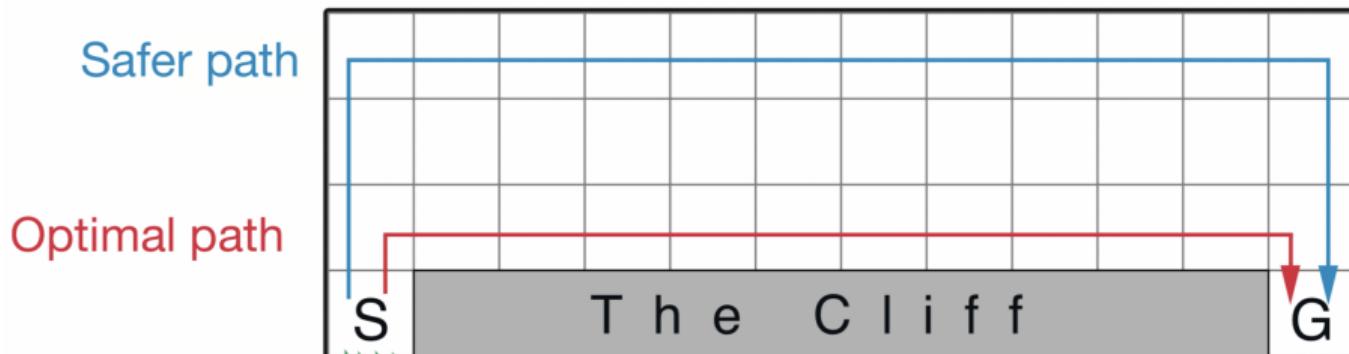


Figure: Illustrations of Cliff Walking

- Undiscounted, **episodic** task
- Actions: up, down, right and left
- Reward of -100 if stepping into **cliff**
- Reward of -1 on other transitions

# Cliff Walking Example (Cont'd)

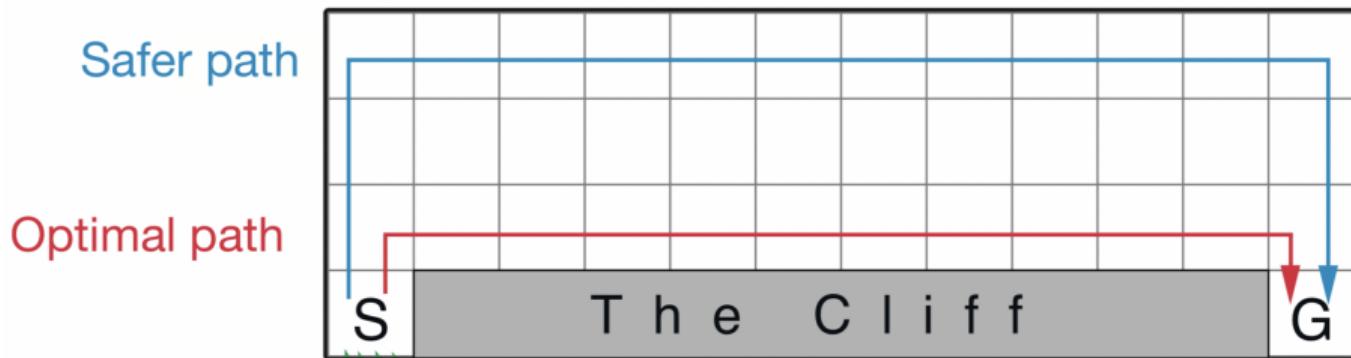


Figure: Illustrations of Cliff Walking

- Q-learning identifies the **optimal** path
- SARSA identifies a **safer** path (the optimal path is not optimal here due to that the  $\epsilon$ -greedy policy, which might force the agent to fall into the cliff when walking along the optimal path, yielding a low value)

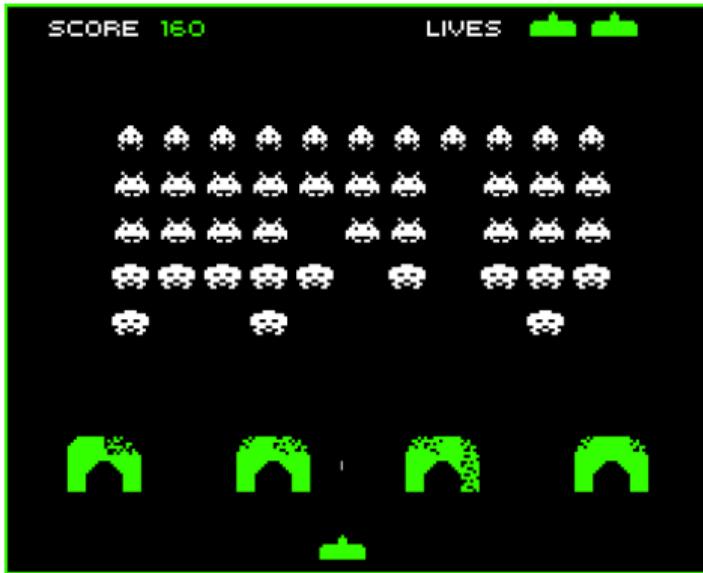
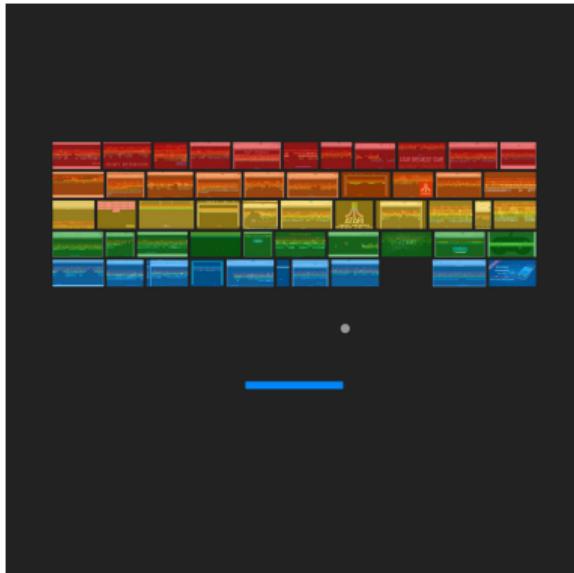
# Limitations of Tabular Methods

---

- So far, we studied reinforcement learning methods using a **tabular representation**
  - Focus on **finite** MDPs
  - Value function represented by a **table**
  - Each state  $s$  has an entry for value  $V(s)$
  - Each state-action pair  $(a, s)$  has an entry for value  $Q(s, a)$
- **Limitations** of tabular methods
  - Cannot handle **large-scale** RL problems or **continuous** state space
  - **Scalability**: storage needed to maintain estimates
  - **Slow learning**: learning the value of each state individually

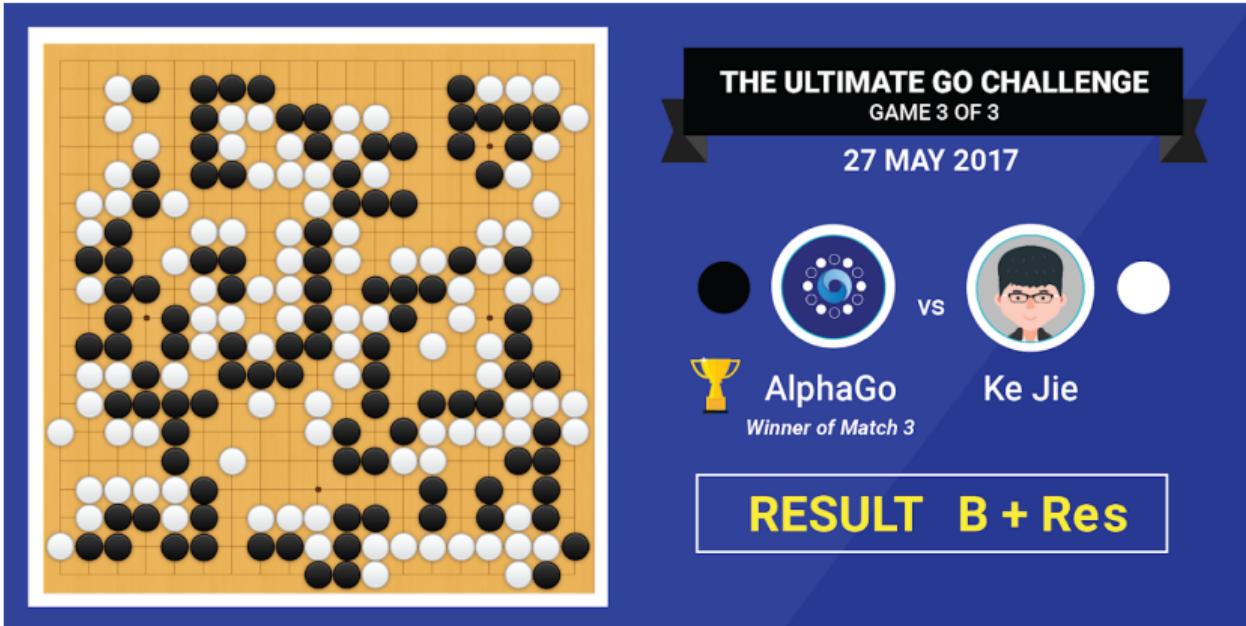
# Large Scale RL Problems (Examples)

---



- Image-valued states (e.g., 210× 160 pixel image frames, 129 colours)

# Large Scale RL Problems (Examples)



- $19 \times 19 = 361$  Go board, each location (empty, black or white)  $\rightarrow 3^{361} \approx 10^{170}$  states

# Lecture Outline

---

1. Tabular Q-Learning
2. Fitted Q-Iteration
3. Case Study I: Deep Q-Network (DQN) in Atari
4. Case Study II: TD Learning in Ridesharing Platforms

# Efficient Control with Function Approximation

---

- **Batch** (offline) setting with pre-collected data  $\{S_t, A_t, R_t, S_{t+1}\}_t$
- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \middle| S_t, A_t \right]$$

- Supervised learning is sample efficient in **batch** settings
- Use supervised learning to learn  $Q^{\pi^{\text{opt}}}$  by solving Bellman optimality equation

# Challenge

---

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \middle| S_t, A_t \right]$$

- Both LHS and RHS involve  $Q^{\pi^{\text{opt}}}$
- A naive approach: minimize the **mean squared Bellman error**

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

This would yield a **biased** estimator!

# The Bellman Error is Not Learnable

---

- For a given random variable  $Z$ ,  $\mathbb{E}Z^2 = (\mathbb{E}Z)^2 + \text{Var}(Z)$
- The mean squared Bellman error can be decomposed into **squared bias** + **variance**

$$\begin{aligned}& \mathbb{E} \left[ \left\{ \mathbf{R}_t + \gamma \max_{\mathbf{a}} \mathbf{Q}(\mathbf{S}_{t+1}, \mathbf{a}) - \mathbf{Q}(\mathbf{S}_t, \mathbf{A}_t) \right\}^2 \middle| \mathbf{A}_t, \mathbf{S}_t \right] \\&= \left[ \mathbb{E} \left\{ \mathbf{R}_t + \gamma \max_{\mathbf{a}} \mathbf{Q}(\mathbf{S}_{t+1}, \mathbf{a}) - \mathbf{Q}(\mathbf{S}_t, \mathbf{A}_t) \middle| \mathbf{A}_t, \mathbf{S}_t \right\} \right]^2 \\&\quad + \text{Var} \left[ \mathbf{R}_t + \gamma \max_{\mathbf{a}} \mathbf{Q}(\mathbf{S}_{t+1}, \mathbf{a}) - \mathbf{Q}(\mathbf{S}_t, \mathbf{A}_t) \middle| \mathbf{A}_t, \mathbf{S}_t \right]\end{aligned}$$

- The second line is **zero** when  $\mathbf{Q} = \mathbf{Q}^{\pi^{\text{opt}}}$
- The third line is **nonzero** for any  $\mathbf{Q}$  and is a function of  $\mathbf{Q}$  as well
- There is no guarantee  $\mathbf{Q}^{\pi^{\text{opt}}}$  is the minimizer

# Fitted Q-Iteration [Riedmiller, 2005]

---

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \mid S_t, A_t \right]$$

Both LHS and RHS involve  $Q^{\pi^{\text{opt}}}$

- **Main idea:** Fix  $Q^{\pi^{\text{opt}}}$  on the RHS

- **Repeat** the following

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  $\tilde{Q} = \hat{Q}$

# Fitted Q-Iteration (Cont'd)

---

- During each iteration, consider the objective function

$$\begin{aligned} & \mathbb{E} \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right]^2 \\ = & \left[ \mathbb{E} \left\{ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right\} \right]^2 \\ + & \text{Var} \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right] \end{aligned}$$

- When  $\tilde{Q}$  is close to  $Q^{\pi^{\text{opt}}}$ , the second line is **small** when  $Q = Q^{\pi^{\text{opt}}}$
- The third line is the same for any  $Q$ , since  $\tilde{Q}$  is fixed

# Fitted Q-Iteration: Algorithm

---

- **Initialization:**  $\hat{Q}, \tilde{Q}$  arbitrary,  $k = 0$
- While ( $k < K$ ) Repeat

Generated data  $\{(S_t, A_t, R_t, S_{t+1})\}$  using policy derived from  $\hat{Q}$  (e.g.,  $\varepsilon$ -greedy)

Compute  $\hat{Q}$  as the argmin of

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

Set  $\tilde{Q} = \hat{Q}$

# Advantages of Fitted Q-Iteration

---

- **Flexibility:** any supervised learning method (e.g., deep learning, boosting, random forest) is applicable to learn the Q-function during each iteration.
- **Efficiency:** borrows the strength of supervised learning for sample-efficient estimation. Allows high-dimensional state information.

# Theoretical Analysis of Fitted Q-Iteration

---

- Let  $\hat{Q}_k$  denote the Q-estimator during the  $k$ th iteration
- Error decomposition: **bias due to initialization + stochastic estimation error**
- The initialization bias  $\rightarrow \mathbf{0}$  as  $k \rightarrow \infty$
- The estimation error  $\rightarrow \mathbf{0}$  when supervised learning provides a **consistent** estimator at each iteration

# Theoretical Analysis of Fitted Q-Iteration (Cont'd)

---

- At the  $k$ th iteration,

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

- Supervised learning target:

$$Q_k(s, a) = \mathbb{E} \left[ R_t + \gamma \max_a \hat{Q}_{k-1}(S_{t+1}, a) \middle| S_t = s, A_t = a \right]$$

# Theoretical Analysis of Fitted Q-Iteration (Cont'd)

- A key inequality

$$\begin{aligned}\sup_{\textcolor{blue}{s}, \textcolor{red}{a}} |\hat{Q}_k(\textcolor{blue}{s}, \textcolor{red}{a}) - Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})| &\leq \sup_{\textcolor{blue}{s}, \textcolor{red}{a}} |\hat{Q}_k(\textcolor{blue}{s}, \textcolor{red}{a}) - Q_k(\textcolor{blue}{s}, \textcolor{red}{a})| \\ &\quad + \gamma \sup_{\textcolor{blue}{s}, \textcolor{red}{a}} |\hat{Q}_{k-1}(\textcolor{blue}{s}, \textcolor{red}{a}) - Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})|\end{aligned}$$

- Iteratively applying the inequality

$$\begin{aligned}\sup_{\textcolor{blue}{s}, \textcolor{red}{a}} |\hat{Q}_k(\textcolor{blue}{s}, \textcolor{red}{a}) - Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})| &\leq \gamma^k \underbrace{\sup_{\textcolor{blue}{s}, \textcolor{red}{a}} |\hat{Q}_0(\textcolor{blue}{s}, \textcolor{red}{a}) - Q^{\pi^{\text{opt}}}(\textcolor{blue}{s}, \textcolor{red}{a})|}_{\text{Initialization Bias}} \\ &\quad + \frac{1}{1 - \gamma} \sup_{\textcolor{blue}{s}, \textcolor{red}{a}} \max_{j=\{1, \dots, k\}} \underbrace{|\hat{Q}_j(\textcolor{blue}{s}, \textcolor{red}{a}) - Q_j(\textcolor{blue}{s}, \textcolor{red}{a})|}_{\text{Estimation Error}}\end{aligned}$$

# Lecture Outline

---

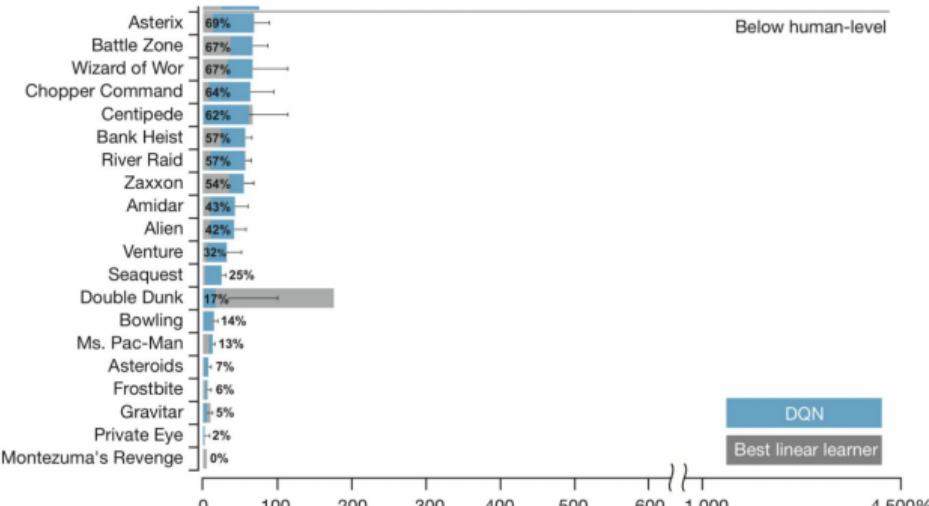
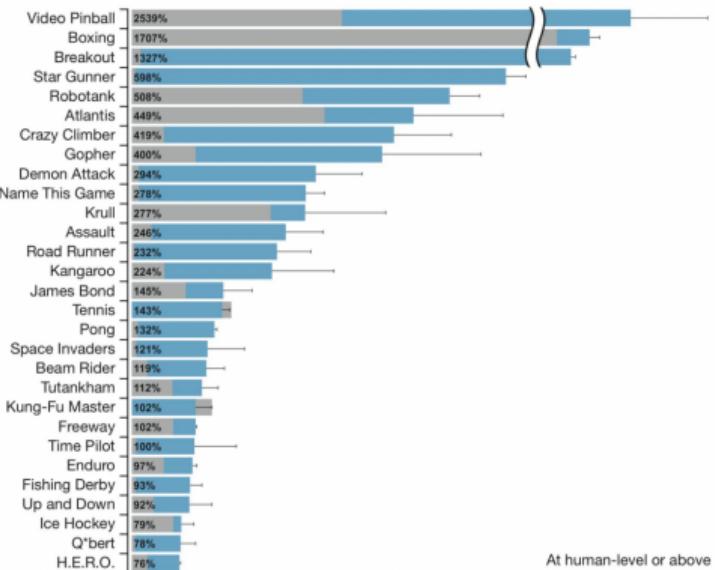
1. Tabular Q-Learning
2. Fitted Q-Iteration
3. Case Study I: Deep Q-Network (DQN) in Atari
4. Case Study II: TD Learning in Ridesharing Platforms

# Deep Q-Network [Mnih et al., 2015]

---

- **Q-learning type method** that uses a neural network Q-function approximator and several tricks to mitigate instability
- Showed superior performance to previously known methods for playing **Atari 2600 games**
- Q-function approximated by a **convolutional neural network**
- Additional tricks: **experience replay, target network**

# DQN: Empirical Results



$$100\% * (\text{DQN score} - \text{random play score}) / (\text{human score} - \text{random play score})$$

# Atari 2600 Observation

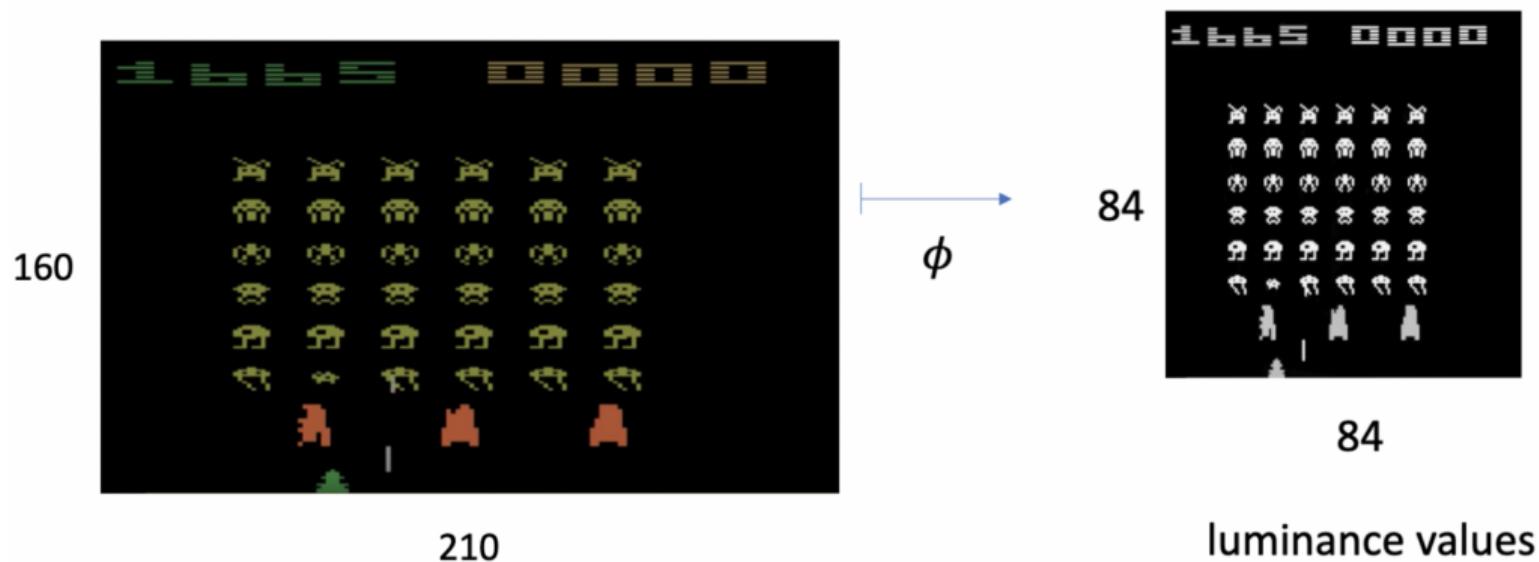
---



- 210×160 pixel image frames
- 129 colours
- 60Hz frame rate
- Non-Markovian

# Input Preprocessing

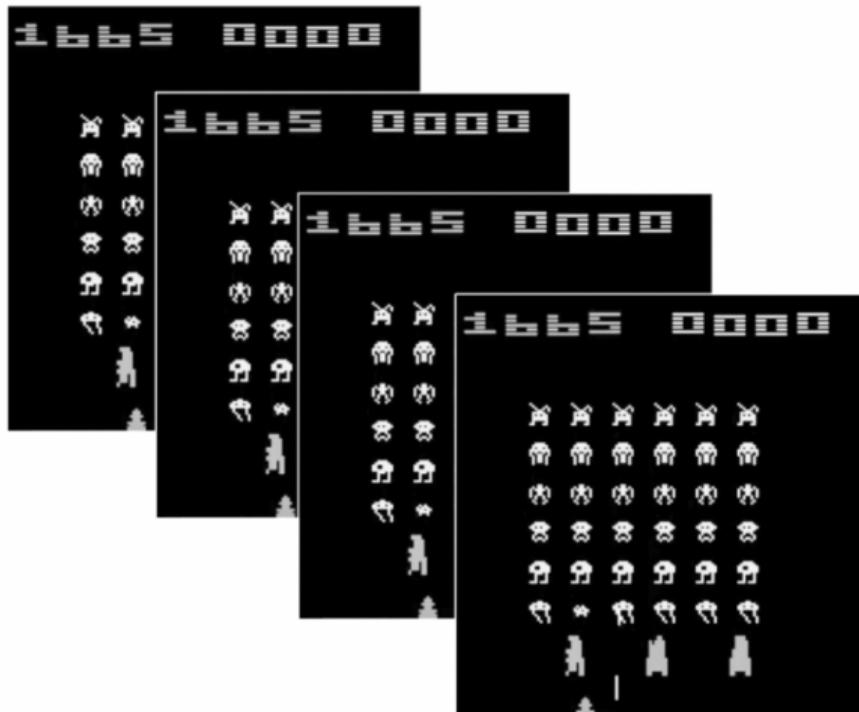
---



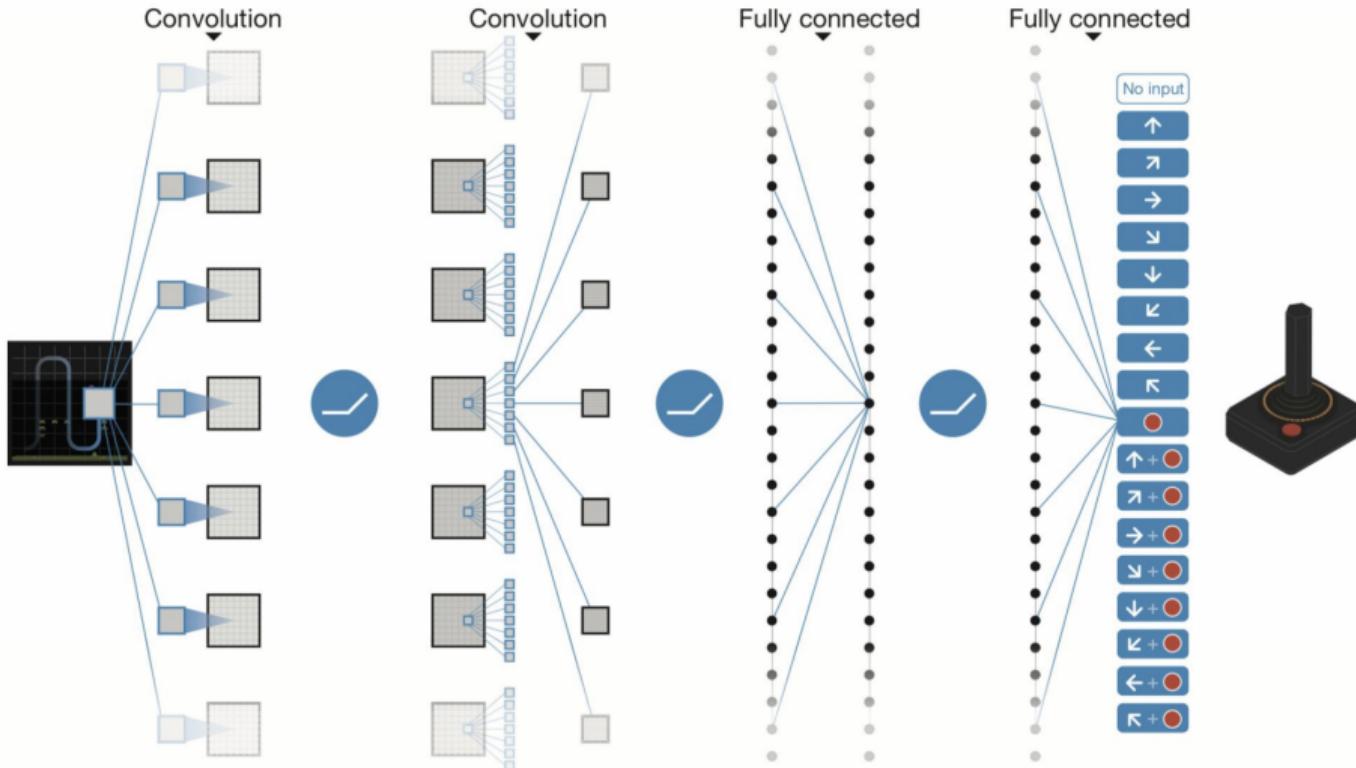
# Mitigating Non-Markovianity by Stacking Frames

---

Input is a stack of 4 most recent frames



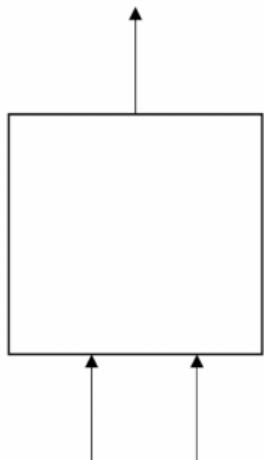
# Action Value Approximator



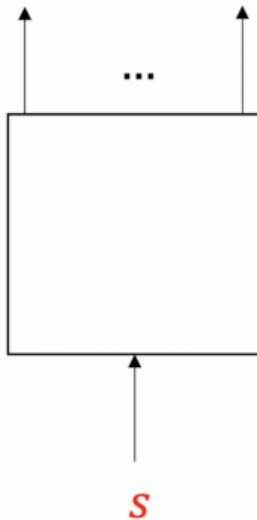
# Action Value Approximator (Cont'd)

---

$$\hat{Q}(s, a; \mathbf{w})$$



$$\hat{Q}(s, a_1; \mathbf{w}) \quad \hat{Q}(s, a_{|A|}; \mathbf{w})$$



# Two Tricks Used in DQN: Experience Replay

---

- **Experience Replay**
    - Store transitions  $(S_t, A_t, R_t, S_{t+1})$  in **replay memory**  $\mathcal{M}$  at time  $t$
    - Sample **minibatch** of transitions  $\{(s_i, a_i, r_i, s_{i+1}) : i \in [n]\}$  from  $n$  and update parameters based on this sub-dataset
    - Differs from tabular Q-learning update
- $$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \max_a [R_t + \gamma Q(S_{t+1}, a) - Q(S_t, A_t)],$$
- where only one tuple is used to update the Q-function.
- Use historical data more **efficiently** to mitigate instability
  - For sufficiently large  $\mathcal{M}$ , the sampled transitions become asymptotically **independent** (correlation decays with time), yielding more accurate estimate

# Recap: Fitted Q-Iteration

---

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \mid S_t, A_t \right]$$

Both LHS and RHS involve  $Q^{\pi^{\text{opt}}}$

- **Main idea:** Fix  $Q^{\pi^{\text{opt}}}$  on the RHS

- **Repeat** the following

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  $\tilde{Q} = \hat{Q}$

# Two Tricks Used in DQN: Target Network

---

- According to Bellman optimality equation

$$\underbrace{Q^{\pi^{\text{opt}}}(S_t, A_t; \theta)}_{\text{Q-network}} = \mathbb{E} \left[ R_t + \gamma \max_a \underbrace{Q^{\pi^{\text{opt}}}(S_{t+1}, a; \theta^*)}_{\text{target-network}} \mid S_t, A_t \right]$$

- Fix  $\theta^*$  in the target network when updating  $\theta$  in the Q-network
- Perform minibatch SGD  $T_{\text{target}}$  steps to update  $\theta$  and set  $\theta^* \leftarrow \theta$
- In Mnih et al. [2015],  $T_{\text{target}} \leftarrow 10000$
- For sufficiently large  $T_{\text{target}}$ , performing minibatch SGD  $T_{\text{target}}$  steps is equivalent to

$$\theta \leftarrow \arg \min_{\theta'} \sum_t \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a; \theta^*) - Q^{\pi^{\text{opt}}}(S_t, A_t; \theta') \right]^2$$

- Share similar spirits as **fitted Q-iteration**

# The Complete Algorithm

---

**Input:** MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ , replay memory  $\mathcal{M}$ , number of iterations  $T$ , minibatch size  $n$ , exploration probability  $\epsilon \in (0, 1)$ , a family of deep Q-networks  $Q_\theta: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , an integer  $T_{\text{target}}$  for updating the target network, and a sequence of stepsizes  $\{\alpha_t\}_{t \geq 0}$ .

Initialize the replay memory  $\mathcal{M}$  to be empty.

Initialize the Q-network with random weights  $\theta$ .

Initialize the weights of the target network with  $\theta^* = \theta$ .

Initialize the initial state  $S_0$ .

**for**  $t = 0, 1, \dots, T$  **do**

With probability  $\epsilon$ , choose  $A_t$  uniformly at random from  $\mathcal{A}$ , and with probability  $1 - \epsilon$ , choose  $A_t$  such that  $Q_\theta(S_t, A_t) = \max_{a \in \mathcal{A}} Q_\theta(S_t, a)$ .

Execute  $A_t$  and observe reward  $R_t$  and the next state  $S_{t+1}$ .

Store transition  $(S_t, A_t, R_t, S_{t+1})$  in  $\mathcal{M}$ .

Experience replay: Sample random minibatch of transitions  $\{(s_i, a_i, r_i, s'_i)\}_{i \in [n]}$  from  $\mathcal{M}$ .

For each  $i \in [n]$ , compute the target  $Y_i = r_i + \gamma \cdot \max_{a \in \mathcal{A}} Q_{\theta^*}(s'_i, a)$ .

Update the Q-network: Perform a gradient descent step

$$\theta \leftarrow \theta - \alpha_t \cdot \frac{1}{n} \sum_{i \in [n]} [Y_i - Q_\theta(s_i, a_i)] \cdot \nabla_\theta Q_\theta(s_i, a_i).$$

Update the target network: Update  $\theta^* \leftarrow \theta$  every  $T_{\text{target}}$  steps.

**end for**

Define policy  $\bar{\pi}$  as the greedy policy with respect to  $Q_\theta$ .

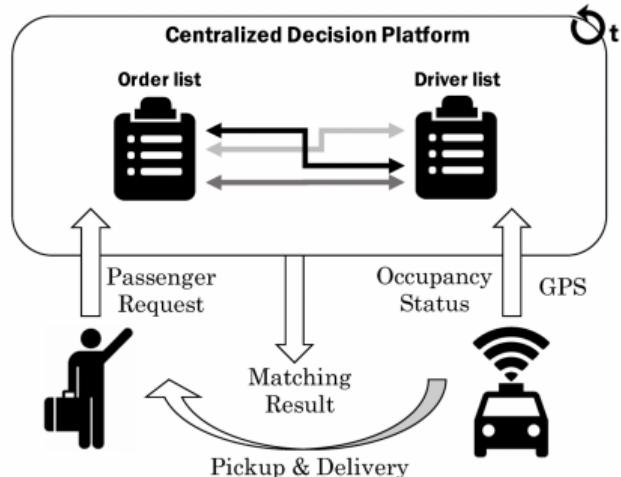
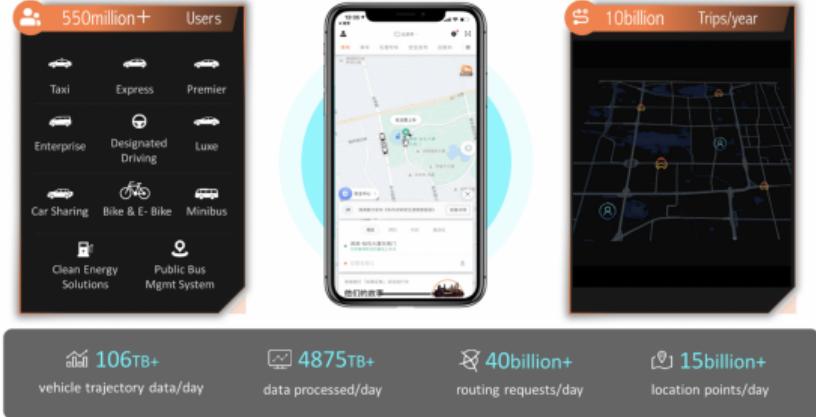
**Output:** Action-value function  $Q_\theta$  and policy  $\bar{\pi}$ .

# Lecture Outline

---

1. Tabular Q-Learning
2. Fitted Q-Iteration
3. Case Study I: Deep Q-Network (DQN) in Atari
4. **Case Study II: TD Learning in Ridesharing Platforms**

# Ridesharing: Order-Dispatching



**Objective:** learn an optimal policy to maximize

- answer rate (proportions of call orders being answered)
- completion rate (proportions of call orders being completed)
- drivers' income

# Order Dispatch Policies

---

- Closest Driver Policy
- MDP Order Dispatch Policy [Xu et al., 2018]
  - **Simple**: no neural networks, no deep learning, use tabular methods
  - **Useful**: performance improvement consistent in all cities, gains in completion rate ranging from 0.5% to 5%, successfully deployed for more than 20 cities
- Some Follow-up Works [Tang et al., 2019, Wan et al., 2021]

# Closest Driver Policy

---

Assign the call order to the closest available driver

$$\arg \min_{\mathbf{a}_{i,j}} \sum_{i=1}^m \sum_{j=1}^n d(i,j) a_{i,j} \quad \text{Minimize driver-passenger total distance}$$

$$s.t. \sum_{i=1}^m a_{i,j} \leq 1, j = 1, \dots, n \quad \text{Order assigned to at most one driver}$$

$$\sum_{j=1}^n a_{i,j} \leq 1, i = 1, \dots, m \quad \text{Driver assigned to at most one order}$$

- $i$  indexes the  $i$ th driver
- $d(i,j)$  = distance between  $i$  and  $j$
- One of the two equalities shall hold
- $j$  indexes the  $j$ th order
- $a_{i,j} = 1 \Leftrightarrow$  order  $j$  is assigned to  $i$

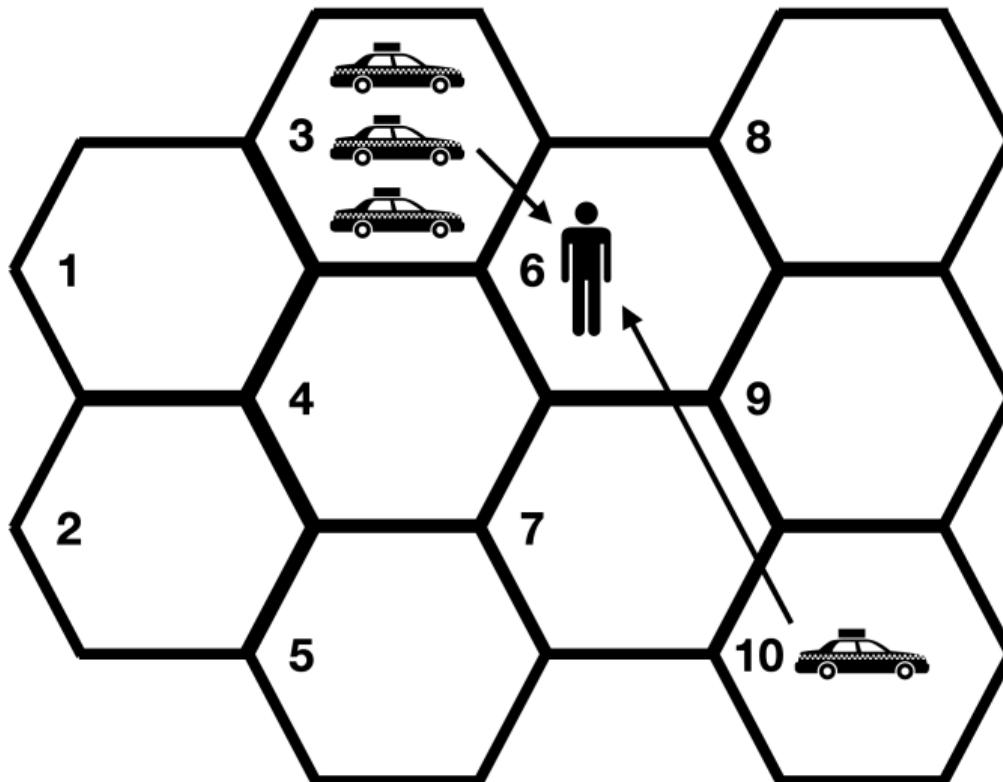
# Closest Driver Policy: Limitations

---

- The company implements the policy every 2 seconds
- **Myopic** policy (e.g., maximize immediate rewards)
- No guarantee it will maximize long-term rewards
- Example given in the next slide

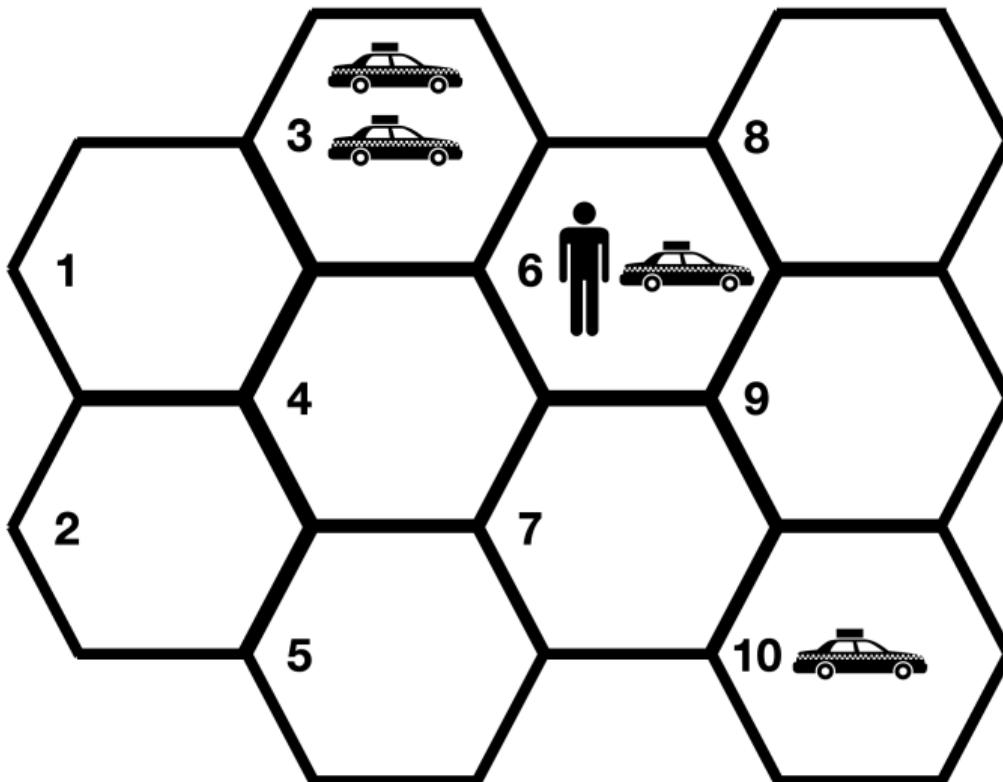
# Illustration of Limitations of Closest Driver Policy

---



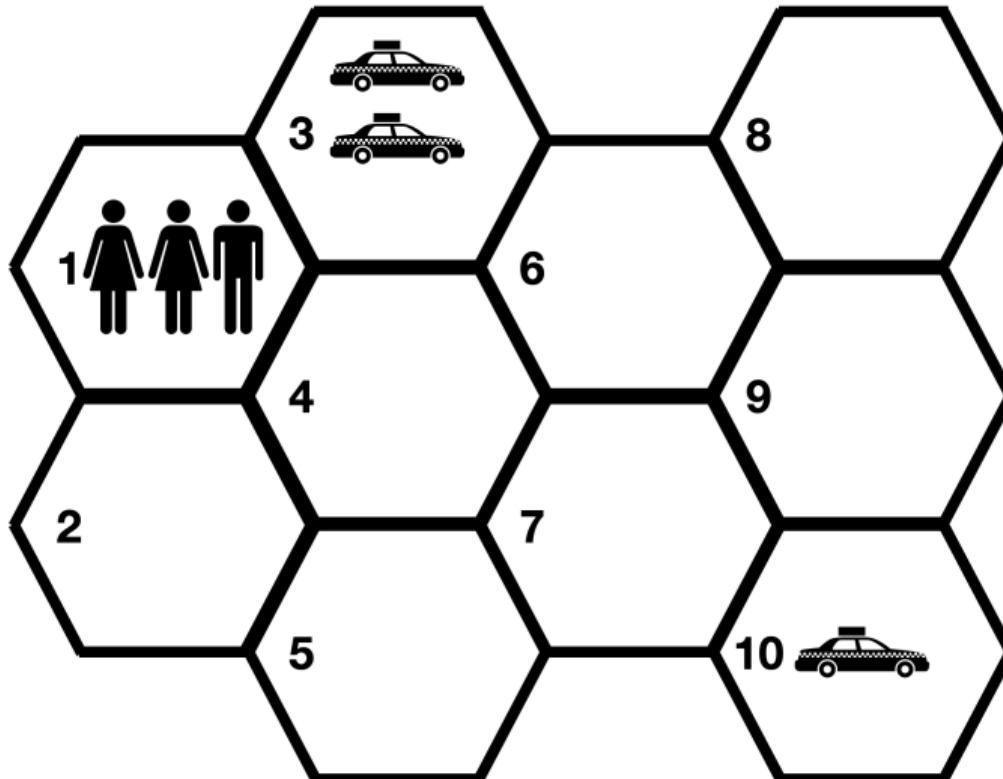
# Adopting the Closest Driver Policy

---



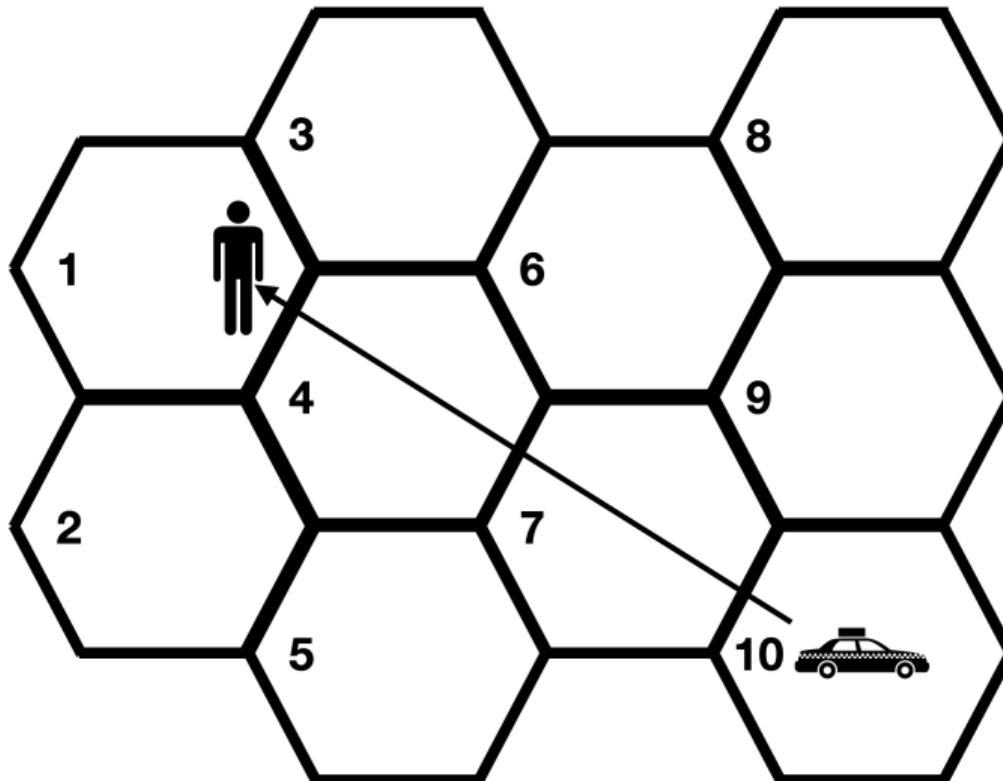
## Some Time Later . . .

---



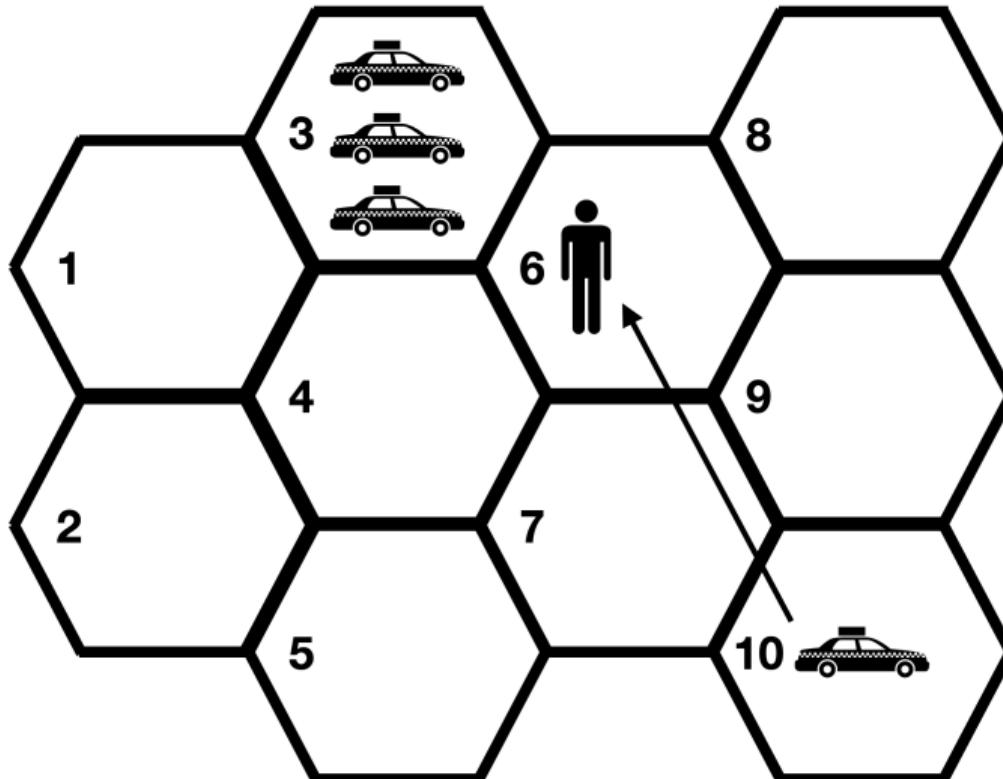
# Miss One Order

---



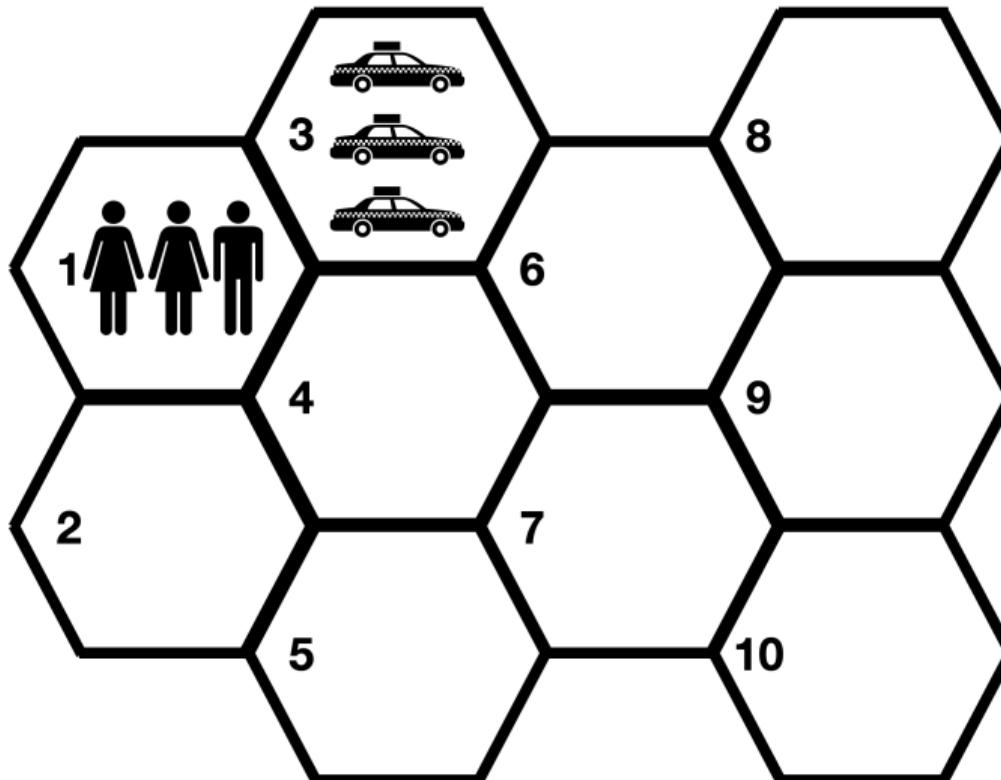
## Consider a Different Action

---



# Able to Match All Orders

---



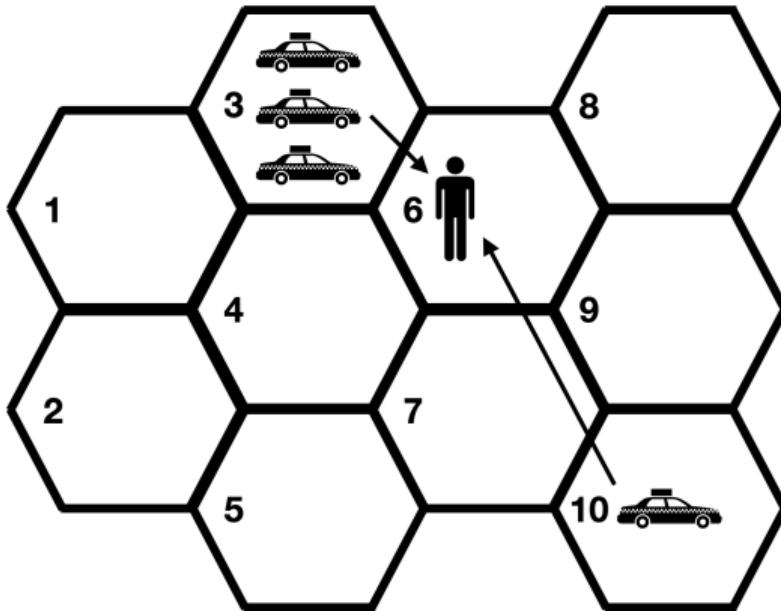
# MDP Order Dispatch Policy

---

- Adopts a **reinforcement learning** framework to optimize long-term rewards
- Delivers **remarkable improvement** on the platform's efficiency
- **Challenges:**
  - Huge state space (e.g., origin/destination of call orders, location of available drivers)
  - Huge action space: number of matchings grows exponentially with number of orders/drivers. With  $n$  orders and  $n$  drivers, number of potential matchings =  $n!$

# Main Idea

---



- Closest driver is myopic because its objective function (e.g., total distance) only considers immediate rewards
- Use an objective function that involves long-term rewards (e.g., value)

# Main Idea (Cont'd)

---

- A learning and planning approach
- **Learning**: policy evaluation based on historical data
- **Planning**: order dispatch by maximizing total value

# An MDP Framework

---

- Model each driver as an agent
- **State**: 2-dim vector (time, location)
- **Action**: two types of actions
  1. Serving action: assign the driver to server an order
  2. Idle action: allows drivers to stay in the same location, to serve an order in the next time
- **Reward**:
  1. an order is completed or not (0/1) (completion rate)
  2. driver's revenue from an order (driver's income)

# An MDP Framework (Cont'd)

---

- **Discounted Factor:**  $\gamma = 0.9$ . An order that lasts for time  $T$  with reward  $R$

$$r = \sum_{t=0}^{T-1} \gamma^t \frac{R}{T}$$

- **Example:**
  - A driver in area **A** receives an order from **B** to **C** at time 00:00
  - The driver arrives **C** at 30min and earns 30£
  - 10min as one time unit,  $\gamma = 0.9$
  - **State transition:**  $(0, A) \rightarrow (3, C)$
  - **Reward:**  $10 + 0.9 \times 10 + 0.9^2 \times 10 = 27.1$

# Learning: Policy Evaluation

---

- Break down historical data into a set of transitions pairs  $(\textcolor{blue}{s}, \textcolor{red}{a}, \textcolor{green}{r}, \textcolor{blue}{s}', \Delta t)$ , where  $\Delta t$  denotes the time of pickup, waiting and deliver process
- TD update rule for the **idle** action

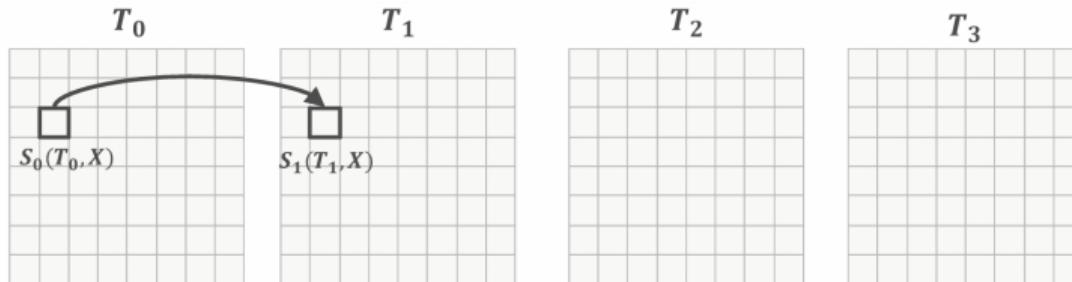
$$V(\textcolor{blue}{s}) \leftarrow V(\textcolor{blue}{s}) + \alpha[\textcolor{green}{0} + \gamma V(\textcolor{blue}{s}') - V(\textcolor{blue}{s})]$$

- TD update rule for the **serving** action

$$V(\textcolor{blue}{s}) \leftarrow V(\textcolor{blue}{s}) + \alpha[\textcolor{green}{r} + \gamma^{\Delta t} V(\textcolor{blue}{s}') - V(\textcolor{blue}{s})]$$

# Policy Evaluation: Example

---



**Idle action:**  $V(S_0) \leftarrow V(S_0) + \alpha(0 + \gamma V(S_1) - V(S_0))$



**Serving action:**  $V(S_0) \leftarrow V(S_0) + \alpha(R_\gamma + \gamma^3 V(S_2) - V(S_0))$

# Policy Evaluation: Pseudocode

---

- **Input:** Collect historical state transitions  $D = \{(s, a, r, s', \Delta t)\}$  where each state is composed of a time and space index
- **Initialize**  $V(s)$  and  $N(s)$  to zero for any  $s$
- **For**  $t = T - 1$  to  $0$  **do**
  - Find a data subset  $D_t$  where the time index of the state is  $t$
  - For** each sample  $(s, a, r, s', \Delta t)$  in  $D_t$  **do**
    - $N(s) \leftarrow N(s) + 1$
    - $V(s) \leftarrow V(s) + N^{-1}(s)[r + \gamma^{\Delta t} V(s') - V(s)]$
  - End For**
- **Return**  $V$

# Planning: Order Dispatch

---

Recall the closest driver policy

$$\arg \min_{\mathbf{a}_{i,j}} \sum_{i=0}^m \sum_{j=1}^n d(i,j) \mathbf{a}_{i,j} \quad \text{Minimize driver-passenger total distance}$$

$$s.t. \quad \sum_{i=1}^m \mathbf{a}_{i,j} \leq 1, \quad j = 1, \dots, n \quad \text{Order assigned to at most one driver}$$

$$\sum_{j=0}^n \mathbf{a}_{i,j} \leq 1, \quad i = 1, \dots, m \quad \text{Driver assigned to at most one order}$$

- $i$  indexes the  $i$ th driver
- $d(i,j)$  = distance between  $i$  and  $j$
- $j$  indexes the  $j$ th order
- $\mathbf{a}_{i,j} = 1 \Leftrightarrow$  order  $j$  is assigned to  $i$

# Planning: Order Dispatch (Cont'd)

---

The MDP order dispatch policy

$$\arg \max_{\mathbf{a}_{i,j}} \sum_{i=0}^m \sum_{j=1}^n \mathbf{A}(i,j) \mathbf{a}_{i,j} \quad \text{Maximize total advantage function}$$

$$s.t. \quad \sum_{i=1}^m \mathbf{a}_{i,j} \leq 1, \quad j = 1, \dots, n \quad \text{Order assigned to at most one driver}$$

$$\sum_{j=0}^n \mathbf{a}_{i,j} \leq 1, \quad i = 1, \dots, m \quad \text{Driver assigned to at most one order}$$

- $i$  indexes the  $i$ th driver
- $\mathbf{A}(i,j) = \text{advantage function}$
- $j$  indexes the  $j$ th order
- $\mathbf{a}_{i,j} = 1 \Leftrightarrow \text{order } j \text{ is assigned to } i$

# Advantage Function Trick

---

- **What** is advantage function?
  - Difference between Q-function and value function.
- **Why** use advantage function trick?
  - Optimize long-term rewards
  - Send drivers in areas with lower values (“cold regions”) to areas with higher values (“hot regions”)

# What is Advantage Function

---

- $A(i, j) = r_{i,j} + \gamma^{\Delta t_{i,j}} V(s'_{i,j}) - V(s_i)$
- $i$  indexes  $i$ th driver,  $j$  indexes  $j$ th order
- $r_{i,j}$ : expected gain for  $i$ th driver to serve  $j$ th order
- $s_i$ : initial state of  $i$ th driver
- $s'_{i,j}$ : state of  $i$ th driver after serving  $j$ th order
  - **Time:**  $s'_{i,j}(t) = s_i(t) + \Delta t_{i,j}$
  - **Location:**  $s'_{i,j}(\ell)$ , the destination of  $j$ th order
- The first two term corresponds to the state-action value function (Q-function) of assigning  $i$ th driver to  $j$ th order

# Why Use Advantage Function Trick

---

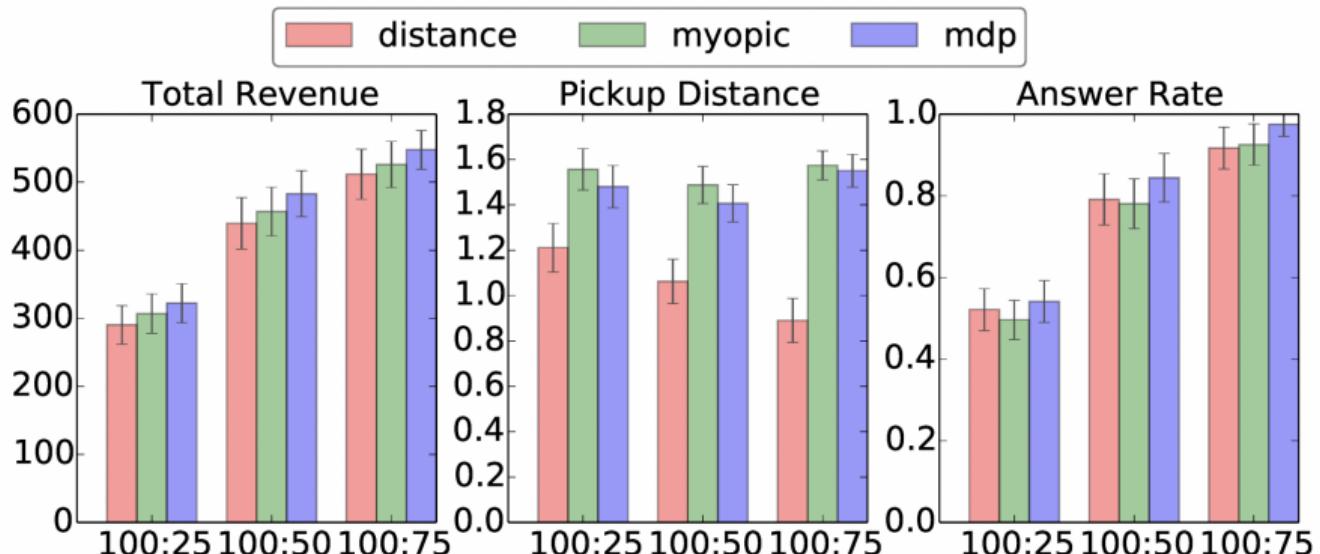
- $A(i, j) = r_{i,j} + \gamma^{\Delta t_{i,j}} V(s'_{i,j}) - V(s_i)$
- **Order Price:** an order with a **high utility** leads to a high advantage
- **Driver's Location:**
  - Value of a driver's current location has a **negative** impact on the advantage
  - When # drivers > # orders (**oversupplied**), drivers in areas with lower values ("cold regions") are more likely to be selected
- **Order's Destination:**
  - Value of an order's destination has a **positive** impact on the advantage
  - When # drivers < # orders (**undersupplied**), orders whose destinations have higher values ("hot regions") are more likely to be selected
- **Pickup Distance:**
  - Contributes to the advantage **implicitly**
  - A larger pickup distance  $\implies$  a larger  $\nabla t_{i,j}$   $\implies$  a lower advantage
  - Considers **immediate reward** as well

# Simulations: Toy Example

---

- A simple map of **9 × 9** spatial grids with **20** time steps
- Orders can only be dispatched to drivers in distance that are no greater than **2**
- Simulate realistic traffic patterns with a morning-peak and a night-peak, centralized on different locations of residential areas and working areas
- Competing methods
  - Distance-based
  - Myopic ( $\gamma = 0$ )

# Toy Example (Cont'd)



**Figure 6: Comparison of distance-based method, myopic method and the proposed MDP method in three metrics on the toy example environment. X-axis stands for the order-driver ratios. Better viewed in color.**

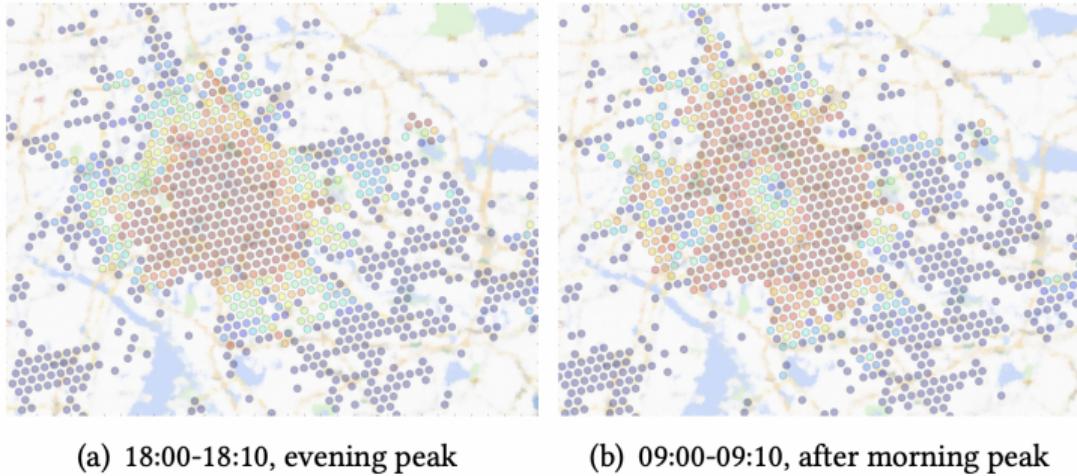
# Real-World Experiment

---

- Performance improvement brought by the MDP method is consistent in all cities
- Gains in global GMV and completion rate ranging from 0.5% to 5%
- Successfully deployed for more than **20** cities
- Serving millions of trips in a daily basis

# Real-World Experiment (Cont'd)

---



**Figure 8: Sampled value function for the same city at different times. Red indicates higher values, blue for lower ones. Better viewed in color.**

# Summary

---

- Tabular Q-Learning
  - Fitted Q-Iteration
  - Deep Q-Network
  - Experience Replay
  - Target Network
- Convolutional Neural Networks
  - Closest Driver Policy
  - MDP Order Dispatch
  - Advantage Function Trick

# References |

---

- Francisco S Melo. Convergence of q-learning: A simple proof. *Institute Of Systems and Robotics, Tech. Rep*, pages 1–4, 2001.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Martin Riedmiller. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European conference on machine learning*, pages 317–328. Springer, 2005.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.

## References II

---

- Xiaocheng Tang, Zhiwei Qin, Fan Zhang, Zhaodong Wang, Zhe Xu, Yintai Ma, Hongtu Zhu, and Jieping Ye. A deep value-network based approach for multi-driver order dispatching. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1780–1790, 2019.
- Runzhe Wan, Sheng Zhang, Chengchun Shi, Shikai Luo, and Rui Song. Pattern transfer learning for reinforcement learning in order dispatching. *arXiv preprint arXiv:2105.13218*, 2021.
- Zhe Xu, Zhixin Li, Qingwen Guan, Dingshui Zhang, Qiang Li, Junxiao Nan, Chunyang Liu, Wei Bian, and Jieping Ye. Large-scale order dispatch in on-demand ride-hailing platforms: A learning and planning approach. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 905–913, 2018.

# Questions