# Optimal Passenger-Seeking Policies on E-hailing Platforms Using Markov Decision Process and Imitation Learning

Zhenyu Shou[a], Xuan Di[a,b,*], Jieping Ye[c], Hongtu Zhu[c], Hua Zhang[d], Robert Hampshire[f,e]

[a]*Department of Civil Engineering and Engineering Mechanics, Columbia University*
[b]*Data Science Institute, Columbia University*
[c]*Didi Chuxing Inc., Beijing, China*
[d]*National Maglev Transportation Engineering R&D Center, Tongji University, Shanghai, China*
[e]*University of Michigan Transportation Research Institute, University of Michigan, Ann Arbor*
[f]*Ford School of Public Policy, University of Michigan, Ann Arbor*

**Abstract**

Vacant taxi drivers' passenger seeking process in a road network generates additional vehicle miles traveled, adding congestion and pollution into the road network and the environment. This paper aims to employ a Markov Decision Process (MDP) to model idle e-hailing drivers' optimal sequential decisions in passenger-seeking. Transportation network companies (TNC) or e-hailing (e.g., Didi, Uber) drivers exhibit different behaviors from traditional taxi drivers because e-hailing drivers do not need to actually search for passengers. Instead, they reposition themselves so that the matching platform can match a passenger. Accordingly, we incorporate e-hailing drivers' new features into our MDP model. The reward function used in the MDP model is uncovered by leveraging an inverse reinforcement learning technique. We then use 44,160 Didi drivers' 3-day trajectories to train the model. To validate the effectiveness of the model, a Monte Carlo simulation is conducted to simulate the performance of drivers under the guidance of the optimal policy, which is then compared with the performance of drivers following one baseline heuristic, namely, the local hotspot strategy. The results show that our model is able to achieve a 17.5% improvement over the local hotspot strategy in terms of the rate of return. The proposed MDP model captures the supply-demand ratio considering the fact that the number of drivers in this study is sufficiently large and thus the number of unmatched orders is assumed to be negligible. To better incorporate the competition among multiple drivers into the model, we have also devised and calibrated a dynamic adjustment strategy of the order matching probability.

*Keywords:* Markov Decision Process (MDP), Imitation Learning, E-hailing

## 1. Motivation

Taxi, complementary to massive transit systems such as bus and subway, provides flexible-route door-to-door mobility service. However, taxi drivers usually have to spend 35-60 percent of their time on cruising to find the next potential passenger (Powell et al., 2011). Such passenger-seeking process not only decreases taxi drivers' income but also generates additional vehicle miles traveled, adding congestion and pollution into the increasingly saturated roads.

Cruising is primarily caused by an imbalance between travel demand and supply. Market regulation (Yang et al., 2002) or taxi fare structure design (Yang et al., 2010a; He et al., 2018; Battifarano and Qian, 2019) were proposed respectively to balance taxi travel demand and supply. A network equilibrium model was developed (Yang and Wong, 1998; Wong and Yang, 1998) to capture the spatial imbalance

---

*Corresponding author. Tel.: +1 212 853 0435;
*Email address:* sharon.di@columbia.edu (Xuan Di)

between travel demand and supply, where a logit-based probability was introduced to describe the meeting between a vacant taxi and a waiting passenger. This model, in which a taxi driver is supposed to minimize the individual search time for the next passenger, is further extended to incorporate congestion effects and customer demand elasticity (Wong et al., 2001), to include the fare structure and fleet size regulation (Yang et al., 2002), to consider multiple user classes, multiple taxi modes, and customer hierarchical modal choice (Wong et al., 2008), and to use a meeting function to describe the search frictions between vacant taxis and waiting passengers (Yang et al., 2010b; Yang and Yang, 2011; Yang et al., 2014). Recently, Di and Ban (2019) proposed a unified equilibrium framework to model the shared mobility in congested road network.

As taxis GPS trajectories become increasingly available, qualitative analysis has been performed to uncover drivers' actual searching strategy. Liu et al. (2010) found that drivers with higher profits prefer to choose routes with higher speed in both operational and idle states. Li et al. (2011) discovered that hunting is a more efficient strategy than waiting by comparing profitable and non-profitable drivers. Several logit-based quantitative models were developed to capture idle drivers' searching behavior (Szeto et al., 2013; Sirisoma et al., 2010; Wong et al., 2014a,b, 2015a,b). The bilateral searching behavior (i.e., taxi searching for customers and customers searching for taxis) was modeled through an absorbing Markov chain approach (Wong et al., 2005). A probabilistic dynamic programming routing model was proposed to capture the taxi driver's routing decisions at intersections (Hu et al., 2012). Furthermore, a two-layer approach, in which the first layer models the driver's pick-up location choice and the second layer accounts for the driver's detailed route choice behavior, was presented (Tang et al., 2016). Recently, Zhang et al. (2019a) proposed an image-based representation of taxi drivers' passenger searching strategies and identified twenty four strategies using a dataset collected in Shenzhen, China.

Upon the understanding of drivers' searching behavior, recommendations can be provided to idle drivers on where to find the next passenger. An accurate prediction of both taxi supply (Phithakkitnukoon et al., 2010) and demand (Moreira-Matias et al., 2012; Markou et al., 2019; Wang et al., 2019; Alemi et al., 2019) as well as travel time (Tan et al., 2018; Zhang et al., 2019b) are stepping stones to these recommendations. The objectives that recommendations aim to achieve include the minimization of waiting time at the recommended location (Hwang et al., 2015) or of the distance between the current location and the recommended location (Powell et al., 2011; Hwang et al., 2015), and the maximization of the expected fare for the next trip (Powell et al., 2011; Hwang et al., 2015), of the probability of finding a passenger (Ge et al., 2010), or of the potential profit of a driver (Qu et al., 2014; Yuan et al., 2011).

The aforementioned studies mainly focused on the recommendation of the cruising routes or next cruising locations at the immediate next step without considering the optimization of long-run payoffs. A recommended customer searching strategy may help a driver to get an order as fast as possible but may not maximize this driver's overall profit in one day. Models which can capture drivers' long-term optimization strategy are needed. In recent years, Markov Decision Process (MDP) becomes increasingly popular in optimizing a single agent's sequential decision-making process given a period of time (Puterman, 1994). Several studies (Rong et al., 2016; Zhou et al., 2018; Verma et al., 2017; Gao et al., 2018; Yu et al., 2019) have employed MDPs to model idle drivers' optimal searching strategy. In an MDP, an idle driver is an agent who makes sequential decisions of where to go to find the next passenger in a stochastic environment. The environment is characterized by a Markov process and transitions from one state to another once an action is specified by the idle driver. The driver aims to select an optimal policy which optimizes her long-term expected reward. Dynamic programming or Q-learning approaches are commonly used to solve an MDP (Sutton and Barto, 1998). Table (1) summarizes the existing studies using MDPs for passenger-seeking optimization. Note that in e-hailing, there is actually no passenger seeking because it is the e-hailing paltform that matches an idle e-hailing driver to a passenger. However, e-hailing drivers still need to *reposition* themselves in order to get better chance of getting matched to a passenger request. In this paper, we will use the terminologies passenger seeking and repositioning interchangably.

Table 1: Existing MDP based models on passenger-seeking strategy

| Reference | Network representation | State space | Action space | Reward | Algorithm |
|---|---|---|---|---|---|
| Rong et al. (2016) | Grid world | (grid id, time, incoming direction) | moving to a neighboring grid or staying in the current grid | Taxi fare | Dynamic programming |
| Verma et al. (2017) | Grid world (static and dynamic zone structure) | (day-of-week, grid id, time-interval) | moving to any chosen grid (proposed an action detection algorithm) | Taxi fare - traveling distance cost - time cost | Q-learning (Monte Carlo) |
| Gao et al. (2018) | Grid world | (grid id, operating status) | driving vacantly to neighboring grids to search, finding a passenger in the current grid, waiting static at the same spot | the ratio of the occupied taxi trip mileage to the previous empty mileage | Q-learning (Temporal Difference) |
| Yu et al. (2019) | Link node | (node id, indicator of the current pickup drop-off cycle) | outgoing links from the current node | taxi fare - operating cost | Value iteration |
| Lin et al. (2018) | Grid world | (grid id, time interval, global state) | moving into a neighboring grid or staying in the current grid | taxi fare - operating cost | Reinforcement Learning (Deep Q learning) |

The existing MDP models were primarily developed for traditional taxi drivers' sequential decision-making where a driver has to see a passenger before a match happens. In other words, an idle driver's searching process ends only when this driver sees a passenger and the passenger accepts the ride (see Figure (1a)). E-hailing applications (such as Didi and Uber), on the other hand, offer an online platform to match a driver with a passenger even when they are not present in the same space at the same time (He and Shen, 2015; Qian and Ukkusuri, 2017). In other words, even when an idle driver sees a passenger waiting on the roadside, as long as the e-hailing platform does not match them, the driver cannot give a ride to the passenger. However, it does not mean e-hailing drivers always stay at the previous drop-off spot and wait for the platform to match. Drivers tend to *reposition* themselves so that the platform can find them a match sooner. As a result, the decision-making process of e-hailing drivers is quite different from the traditional taxi drivers in the following aspects:

1. An e-hailing driver may receive a matched order before she drops off the previous passenger, thus there is no passenger seeking (see Figure (1b)).

2. Different from traditional taxi that a driver has to see a passenger to find a match, e-hailing platforms very likely find a match even when the driver and the passenger are spatially far from each other. In other words, a driver's search process may end before a passenger is picked up (see Figure (1c)).
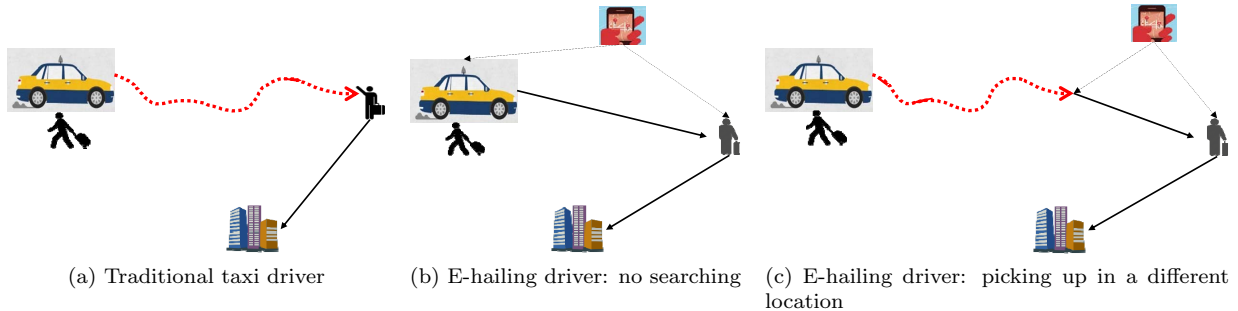
(a) Traditional taxi driver     (b) E-hailing driver: no searching     (c) E-hailing driver: picking up in a different location

Figure 1: Sequential decision-making processes for traditional and e-hailing taxi drivers

Because of the inherent differences in drivers' decision-making, this paper aims to develop an MDP to model e-hailing drivers' sequential decision-making in searching for the next passenger. 44,160 Didi drivers' 3-day GPS trajectories are used to calibrate and validate our model. Previously, there is research using Didi's data for the study of large-scale fleet management (Lin et al., 2018) and large-scale order dispatch (Xu et al., 2018) in e-hailing platforms.

The major contributions of this paper are as follows: (1) In stead of following the literature where a known reward function is given based on some prior knowledge or assumptions (Rong et al., 2016; Verma et al., 2017; Lin et al., 2018; Yu et al., 2019), this work unveils the underlying reward function of the overall e-hailing driver population and crafts a novel reward function which explains the behaviors of drivers with a relatively small radius of gyration and thus paves the way for future research on discovering the underlying reward mechanism in a complex and dynamic e-hailing market. With the incomplete and noisy observed policy, this work first extracts the underlying reward function and then solves an MDP to derive the optimal policy which completes and corrects the observed policy. (2) To the best of our knowledge, this is the first study using large amounts data to devise and calibrate a dynamic adjustment strategy of the order matching probability to address the competition among multiple drivers. The strategy essentially attenuates the order matching probability in an exponential manner for subsequent drivers to be guided into a grid when some drivers have already entered the grid. The strategy is further verified to be efficient in providing different recommendations for multiple drivers.

The remainder of the paper is organized as follows. Section 2 introduces our modified MDP model and details definitions of states, actions, and state transitions and the process of extracting parameters from the data. Section 3 presents the proposed dynamic adjustment strategy of the order matching probability and details the calibration process. Section 4 introduces the data we used in this research and presents the results, including the derived optimal policy and the Monte Carlo simulation. Section 5 concludes the paper and provides some future research directions.

## 2. Markov Decision Process (MDP) for a single agent

### 2.1. Preliminaries

An MDP is specified by a tuple $(S, A, R, P, s_0)$, where $S$ denotes the state space, $A$ stands for the allowable actions, $R$ collects rewards, $P$ defines a state transition matrix, and $s_0$ is the starting state. Given a state $s_t = s \in S$ and a specified action $a_t = a \in A$ at time $t$, the probability of reaching state $s'$ at time $t + 1$ is determined by the probability transition matrix $P(s, a, s')$, which is defined as

$$P(s, a, s') \equiv Pr(s_{t+1} = s' | s_t = s, a_t = a) \tag{2.1}$$

From the initial state $s_0$, the process proceeds repeatedly by following the dynamics of the environment defined by the Equation (2.1) until a terminal state (i.e., either the current time exceeds the terminal time or the current state is an absorbing state) is reached. An MDP satisfies the Markov property which essentially says that the future process is independent on the past given the present, i.e.,

$$Pr(s_{t+1} = s' | s_0, a_0, \cdots, s_{t-1}, a_{t-1}, s_t, a_t) = Pr(s_{t+1} = s' | s_t = s, a_t = a). \tag{2.2}$$

There are two types of value functions in MDPs, namely, a state value $V(s)$ and a state-action value $Q(s, a)$. The actions that an agent will take form a policy $\pi$, which is a mapping from a state $s$ and an

4

action $a$ to the probability $\pi(a|s)$ of taking action $a$ at state $s$. Then the value function of a state $s$ by following the policy $\pi$, denoted as $V_\pi(s)$, can be taken as the expectation of the future rewards, i.e.,

$$V_\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{K} \gamma^k r_{t+k+1} | s_t = s \right] \tag{2.3}$$

where $\gamma$ is a discount factor. The state-action value of taking action $a$ at state $s$ by following policy $\pi$ is

$$Q_\pi(s,a) = \mathbb{E}_\pi \left[ \sum_{k=0}^{K} \gamma^k r_{t+k+1} | s_t = s, a_t = a \right] \tag{2.4}$$

The value function $V_\pi(s)$ is actually a weighted average of the state-action value $Q_\pi(s,a)$, i.e.,

$$V_\pi(s) = \sum_{j=1}^{J} \pi(a = a_j | s) Q_\pi(s, a = a_j) \tag{2.5}$$

where $\pi(a = a_j | s)$ is again the probability of taking action $a_j$ at state $s$ according to policy $\pi$, and $J$ is the total number of actions that are allowed to be taken in state $s$.

Several algorithms have been developed to solve the MDP, i.e., to derive the optimal policy, and the corresponding value functions, such as the dynamic programming method and the Q-learning approach (Sutton and Barto, 1998). The dynamic programming algorithm is used in this work and will be explained later in Section (2.2.4). Furthermore, these two types of value functions at optimality are related by the following mechanism

$$V(s) = max_a Q(s,a) \tag{2.6}$$

The rationale underlying the relationship at optimality is simply to choose a policy which maximizes the value function expressed in Equation (2.5). For example, when an agent is at state $s$, the optimal policy simply suggests the agent to take an action $a$ with the largest state-action value, i.e., $\pi(\arg\max_a Q(s,a)|s) = 1$ (i.e., the probability of taking action $\arg\max_a Q(s,a)$ in state $s$ is 1). Accordingly, the state-action value at optimality can be written as

$$Q(s,a) = \sum_{s'} Pr(s'|s,a)(V(s') + r(s,a,s')) \tag{2.7}$$

where $Pr(s'|s,a)$ is the probability of landing in state $s'$ after taking action $a$ in state $s$, and $r(s,a,s')$ is the reward for choosing action $a$ at state $s$ and landing in state $s'$.

*2.2. MDP for e-hailing drivers*

In this section, we will develop an MDP model for e-hailing drivers' stochastic passenger seeking process. Notations which will be used in the subsequent analysis are listed in Table (2).

Table 2: Notations

| Variable | Explanation |
|---|---|
| $l$ | Index of the current grid |
| $t$ | Current time |
| $I$ | Indicator, denoting whether the driver has been matched to a request before the next drop-off |
| $s$ | State, $s = (l, t, I)$ |
| $S$ | State space, a collection of all states |
| $a$ | Action |
| $A$ | Action space |
| $t_{seek}(l_a)$ | Time spent on seeking for a passenger in grid $l_a$ |
| $t_{drive}(l, k)$ | Time spent on moving from grid $l$ to grid $k$ |
| $d_{seek}(l_a)$ | Distance traveled when seeking for a passenger in grid $l$ |
| $d_{drive}(l, k)$ | Distance traveled for moving from grid $l$ to grid $k$ |
| $p_{order\_match}(l_a)$ | The probability that the driver can be matched to a request during cruising in grid $l_a$ |
| $p_{pickup}(l_a, l'')$ | The probability of picking up a passenger in grid $l''$ when the request from the passenger was matched to the driver in grid $l_a$ |
| $p_{dest}(l'', l''')$ | The probability of dropping off a passenger in grid $l'''$ when the passenger was picked up in grid $l''$ |
| $p_{match}(l'', l''')$ | The probability of receiving a new request before the driver finishing her current order at grid $l'''$ |
| $f(l'', l''')$ | The average taxi fare from grid $l''$ to grid $l'''$ |
| $\alpha$ | Coefficient of fuel consumption and other operating costs per unit distance |

*2.2.1. States*

In our MDP model, the state $s = (l, t, I)$ consists of three components, namely, a grid index $l \in L$, current time $t \in T$, and an indicator $I \in \{0, 1\}$. Note that a hexagonal grid world setting with 6,421 grids is adopted in this research and will be explained later. Considering the fact that an e-hailing driver may receive a request before she drops off the previous passenger, we have therefore added an indicator into the state. The indicator denotes whether the driver has been matched to a request before she arrives at the current state. Accordingly, states with indicator 0 are decision-making states in which the driver needs to spend time on seeking the next passenger, and states with indicator 1 are non-decision-making states. For example, $(1, 2, 1)$ is a non-decision-making state which says that the driver is in grid 1 when $t = 2$ and the driver has already been matched to a request so she will not spend time on seeking at the current state.

*2.2.2. Actions*

In decision-making states, the driver has to choose one from eight allowable actions, denoted as $A$. In non-decision-making states, the driver will not take any action but drive to pick up the next passenger and transport the passenger to the destination. Among the allowable action space $A$, each of the first six actions is to transit from the current grid to one of the six neighbor grids. Note that some of the six neighboring grids may be non-reachable, we thus add a large penalty, i.e., a large distance, to the transition from a grid to a non-reachable neighboring grid to prevent the agent from taking the action which leads the agent to the non-reachable neighboring grid. The seventh action is to stay and cruise around within the current grid. The last action is to wait in the current grid. We stress that the last two actions are essentially different because from the data we have observed that some drivers will just wait near the previous drop-off spot, especially when they are around downtown or transportation terminals while some drivers usually cruise within the current grid after completing a ride. In addition, the fuel cost associated with waiting can be neglected while that of staying can be substantial because the driver keeps cruising around during his/her staying in the current grid. Furthermore, drivers can take a rest and refresh their minds during waiting and hence their driving strategy can be more efficient for future trips. These arguments, however, do not necessarily suggest that the driver should always choose waiting rather than staying. Actually, drivers have to cruise around to get closer to the potential requests under

certain circumstances.

### 2.2.3. State transition

After completing a ride, there are two possible scenarios according to two different values of the indicator. If the indicator is 0, the driver needs to specify an action, i.e. where to find the next passenger, and then moves into the grid along the direction defined by the action and spends some amount of time seeking for the next passenger in the new grid. There are two possible outcomes associated with this passenger seeking process. Either the driver confirms a request and ends up arriving at a different grid by following the passenger's travel plan or the driver fails to find a request and stays in the current grid. The reward is usually positive for the former while negative for the latter due to the fuel consumption and other operating costs. If the indicator is 1, the driver drives to the pick-up spot and then transports the passenger to the destination without any passenger seeking involved.
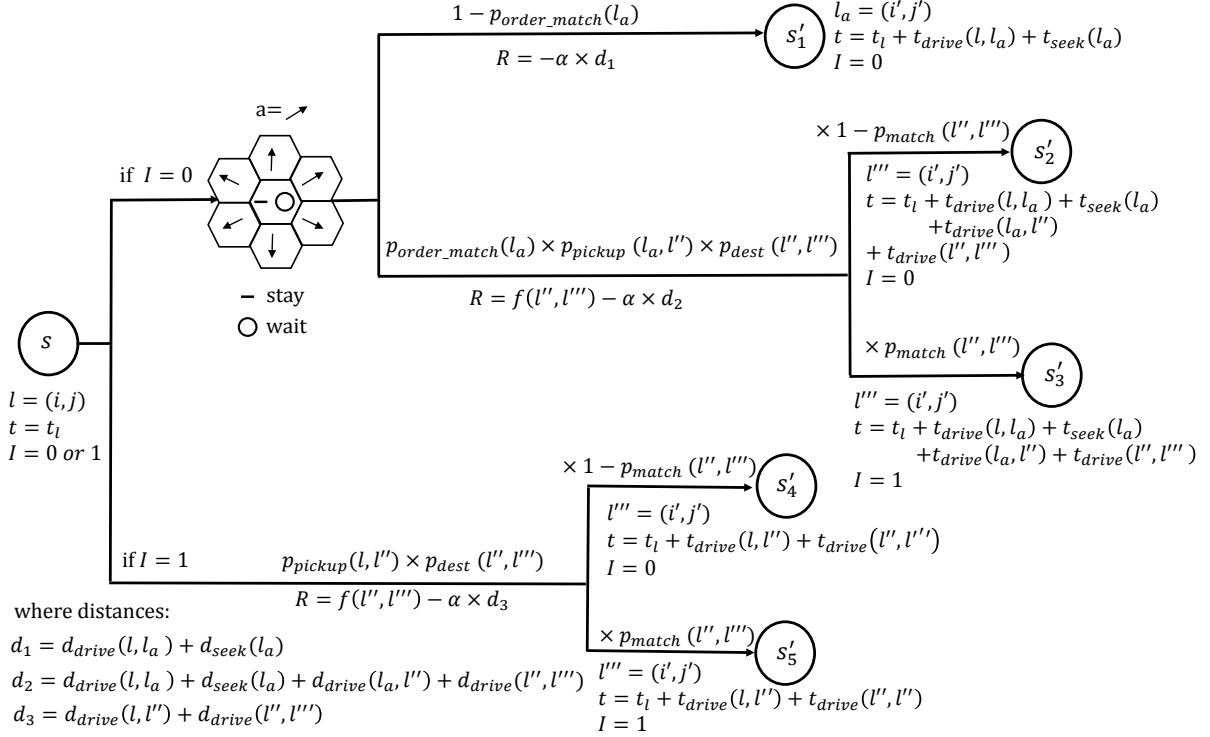


Figure 2: MDP state transition

Figure (2) illustrates the aforementioned state transition process. The driver currently stays at state $s = (l, t, I)$.

If $I = 0$, the driver specifies an action $a$, which is assumably taken as going northeast in the demonstration. Then the driver moves into the grid $l_a$ along the direction defined by the action $a$ and thereafter spends some amount of time $t_{seek}(l_a)$ on seeking the next passenger. There are two possible outcomes of this passenger seeking process.

The first possibility is that the driver fails to get any request in grid $l_a$ after $t_{seek}(l_a)$. In this case, the state of the driver will be $s'_1 = (l', t + t_{drive}(l, l_a) + t_{seek}(l_a), 0)$. The reward for this passenger seeking process is $R = -\alpha(d_{drive}(l, l_a) + d_{seek}(l_a))$, which is negative. Let $P_{order\_match}(l_a)$ denote the probability that the driver will receive at least one request in grid $l_a$. Then the probability of the occurrence of this outcome is $1 - P_{order\_match}(l_a)$. In other words, with probability $1 - P_{order\_match}(l_a)$, the driver will end up in state $s'_1 = (l_a, t + t_{drive}(l, l_a) + t_{seek}(l_a), 0)$.

The second possibility is that the driver confirms one request during the cruising process in $l_a$. The probability of the occurrence of this outcome is $P_{order\_match}(l_a)$. We let $P_{pickup}(l_a, l'')$ denote the probability of confirming a request in grid $l_a$ and picking up the passenger in grid $l''$. Once the passenger is on board, the driver will directly move to the destination $l'''$, which only depends on the passenger's travel plan. We let $P_{dest}(l'', l''')$ denote the probability of picking up a passenger in grid $l''$ and dropping off the passenger in grid $l'''$. After dropping off the passenger, the driver will end

7

up in grid $l'''$ at time $t + t_{drive}(l, l_a) + t_{seek}(l_a) + t_{drive}(l_a, l'') + t_{drive}(l'', l''')$ and earn a reward of $r(l, l''') = f(l'', l''') - \alpha(d_{drive}(l, l_a) + d_{seek}(l_a) + d_{drive}(l_a, l'') + d_{drive}(l'', l'''))$. Hence, the probability of the driver receives a request in grid $l_a$, picks up the passenger in grid $l''$, and transports the passenger to grid $l'''$ is $P_{order\_match}(l_a) \times P_{pickup}(l_a, l'') \times P_{dest}(l'', l''')$. Notice that for a driver, during her trip from the passenger's origin $l''$ to the passenger's destination $l'''$, there is a probability at which the driver will confirm a request before she drops off the passenger. Let $P_{match}(l'', l''')$ denote the probability of receiving a request before the driver reaches grid $l'''$. Then we can conclude that with probability $P_{order\_match}(l_a) \times P_{pickup}(l_a, l'') \times P_{dest}(l'', l''') \times (1 - P_{match}(l'', l'''))$, the driver will end up in state $s_2' = (l''', t + t_{drive}(l, l_a) + t_{seek}(l_a) + t_{drive}(l_a, l'') + t_{drive}(l'', l'''), 0)$, and with probability $P_{order\_match}(l_a) \times P_{pickup}(l_a, l'') \times P_{dest}(l'', l''') \times P_{match}(l'', l''')$, the driver will end up in state $s_3' = (l''', t + t_{drive}(l, l_a) + t_{seek}(l_a) + t_{drive}(l_a, l'') + t_{drive}(l'', l'''), 1)$.

If $I = 1$, the driver will not need to specify any action and will directly drive to the pick-up spot of the next passenger and then transport the passenger to the destination. Again, during her trip to the passenger's destination, there is a probability, denoted as $P_{match}(l'', l''')$, at which the driver will receive a request before she drops off the passenger. As illustrated in Figure (2), with probability $P_{pickup}(l, l'') \times P_{dest}(l'', l''') \times (1 - P_{match}(l'', l'''))$, the driver will end up in state $s_4' = (l''', t + t_{drive}(l, l'') + t_{drive}(l'', l'''), 0)$; with probability $P_{pickup}(l, l'') \times P_{dest}(l'', l''') \times P_{match}(l'', l''')$, the driver will end up in state $s_5' = (l''', t + t_{drive}(l, l'') + t_{drive}(l'', l'''), 1)$.

In both scenarios, namely, either $I = 0$ or $I = 1$, the driver will thereafter start the whole process from $s'$ again until the current time exceeds the time interval, i.e., a terminal state has been reached.

*2.2.4. Solving MDP*

The objective of the MDP model is to maximize the total expected revenue of a driver. Considering the fact that in a time interval, a driver can finish a finite number of pick-up and drop-off cycles, indicating that the MDP model is finite-horizon. When current time of the driver has reached the end of the time interval, no more actions can be taken and no more rewards can be earned. Suppose a driver is currently at state $s = (l, t, I)$. If $I = 0$, meaning that the driver is at a decision-making state, the maximum expected revenue that a driver can earn by starting from $s$ and specifying an action $a$ is

$$
\begin{aligned}
Q(s, a) \quad = \quad & (1 - P_{order\_match}(l_a)) \times [-\alpha(d_{drive}(l, l_a) + d_{seek}(l_a)) + V^*(s_1')] \\
& + \sum_{l'' \in L} \sum_{l''' \in L} P_{order\_match}(l_a) \times P_{pickup}(l_a, l'') \times P_{dest}(l'', l''') \\
& \times [f(l'', l''') - \alpha(d_{drive}(l, l_a) + d_{seek}(l_a) + d_{drive}(l_a, l'') + d_{drive}(l'', l''')) \\
& + (1 - P_{match}(l''')) \times V^*(s_2') + P_{match}(l''') \times V^*(s_3')]
\end{aligned}
\tag{2.8}
$$

where $l_a = l_a(s, a)$, meaning that the grid $l_a$ in which an e-hailing driver will be cruising is dependent on the current state $s$, actually through the grid index $l$ of $s$, and the specified action $a$, $s_1' = (l_a, t + t_{drive}(l, l_a) + t_{seek}(l_a), 0)$, $s_2' = (l''', t + t_{drive}(l, l_a) + t_{seek}(l_a) + t_{drive}(l_a, l'') + t_{drive}(l'', l'''), 0)$, $s_3' = (l''', t + t_{drive}(l, l_a) + t_{seek}(l_a) + t_{drive}(l_a, l'') + t_{drive}(l'', l'''), 1)$, and $V^*(s_1')$, $V^*(s_2')$, and $V^*(s_3')$ stand for the maximum expected revenue that a driver can earn by reaching state $s_1'$, $s_2'$, and $s_3'$, respectively. If $I = 1$, meaning that the driver is at a non-decision-making state, the driver will not specify any action, and the expected revenue that the driver can earn is

$$
\begin{aligned}
Q(s, .) \quad = \quad & \sum_{l'' \in L} \sum_{l''' \in L} P_{pickup}(l, l'') \times P_{dest}(l'', l''') \\
& \times [f(l'', l''') - \alpha(d_{drive}(l, l'') + d_{drive}(l'', l''')) \\
& + (1 - P_{match}(l'', l''')) \times V^*(s_4') + P_{match}(l'', l''') \times V^*(s_5')]
\end{aligned}
\tag{2.9}
$$

where $s_4' = (l''', t + t_{drive}(l, l'') + t_{drive}(l'', l'''), 0)$, $s_5' = (l''', t + t_{drive}(l, l'') + t_{drive}(l'', l'''), 1)$, and $V^*(s_4')$ and $V^*(s_5')$ stand for the maximum expected revenue that a driver can earn by reaching state $s_4'$ and $s_5'$, respectively.

Then the optimal policy for a driver to follow at a decision-making state $s$ is

$$
\pi(s) = argmax_a[Q(s, a)],
\tag{2.10}
$$

and the maximum expected revenue that a driver can earn by reaching state $s$ is

$$V^*(s) = \begin{cases} max_a Q(s,a) & \text{if } s \text{ is a decision-making state} \\ Q(s,.) & \text{if } s \text{ is a non-decision-making state} \end{cases} \tag{2.11}$$

The policy in Equation (2.10) is deterministic, meaning that the driver can only take one action at the current decision-making state $s$ if she follows the policy. Actually here we slightly abuse the notation. The policy is supposed to be $\pi(argmax_a[Q(s,a)]|s) = 1$, i.e., the probability of taking action $argmax_a[Q(s,a)]$ at state $s$ is 1. It is equivalent to say that at state $s$, the action to take is $argmax_a[Q(s,a)]$, and thus we write the policy at state $s$ as Equation (2.10). A deterministic policy defines a one-to-one mapping from a state to an action. The deterministic policy works when there is only one driver who learns the optimal policy and follows the policy. Otherwise, there might be excess taxi supply at some areas, resulting in a localized competition among taxis. A circulating mechanism was employed to tackle this overload problem (Ge et al., 2010). A multi-agent reinforcement learning approach (Lin et al., 2018) was proposed to consider the competition among drivers. In this research, we use a dynamic adjustment strategy to update the order matching probability when multiple idling e-hailing drivers are guided into the same grid. The proposed dynamic adjustment strategy will be introduced in Section (3).

To efficiently solve the MDP, i.e., to derive an optimal policy, a dynamic programming approach is employed (Bertsekas, 2000; Sutton and Barto, 1998). The basic idea of the dynamic programming algorithm is to divide the overall problem into subproblems and hence to make use of the results of the subproblems to solve the overall problem. An important advantage of the dynamic programming algorithm is that it caches results of all subproblems and thus it is guaranteed that the same subproblem is only solved once.

Now we elucidate how we apply the dynamic programming algorithm to solve the MDP. The goal is to solve the optimal value for all states $s = (l,t,I)$, where $l \in L$, $t \in \{0,1,2,\cdot,180\}$, and $I \in \{0,1\}$. There are in total $6,421 \times 181 \times 2 = 2,324,402$ states, and half of them are decision-making states. We define one subproblem as solving the optimal value for one state and thus we have in total $2,324,402$ subproblems. Noticing that at the final time step, i.e., $t = T = 180$, the maximum expected reward that a driver can earn is obviously zero, we thus have $V^*(s) = 0$ for all states $s$ where $t = T$. For any state $s$ with $t < T$ and a chosen action $a$, the calculation of the state-action value $Q(s,a)$ depends on the value of some future states, i.e., $s'_1$, $s'_2$, and $s'_3$ in Equation (2.8). In other words, the subproblem, i.e., solving the optimal value for state $s$, depends on some subproblems, i.e., solving the optimal value of some future states, e.g., $s'_1$, $s'_2$, and $s'_3$. For a future state $s'$, there might be several states $s$ from which the agent will reach the future state $s'$, indicating the calculation of the optimal value of all these states $s$ requires the calculation of the optimal value of the future state $s'$, resulting in calculating the optimal value of the same state $s'$ multiple times and thus wasting computation power. To avoid the repeated calculation of the optimal value for the same state, we adopt the dynamic programming algorithm. Since the optimal values for all states with $t = T$ are known and the optimal value of a state $s$ depends on the optimal value of some future states, we solve the optimal value of states backwards in time and simply store the solved optimal values in a hash table. Then for a state $s$ and a chosen action $a$, we simply read the optimal values of future states $s'_1$, $s'_2$, and $s'_3$ from the hash table and use Equation (2.8) to calculate the state-action value $Q(s,a)$, based on which the optimal value of the state $s$ can be derived from Equation (2.11). The pesudo code is in Algorithm (1).

**Algorithm 1** Dynamic programming algorithm
___
1: Input: $L$, $T = 180$, $A$, $P_{order\_match}$, $P_{pickup}$, $P_{dest}$, $P_{match}$, fare $f$, $d_{drive}$, $d_{seek}$, $t_{drive}$, $t_{seek}$
2: **Initialize:** a hash table $V^*$ to store optimal values
3: $V^*(s) = 0$ for all states $s$ where $t = T$
4: **for** $t = T - 1$ **to** 1 **do**
5:     **for** $l \in L$ **do**
6:         Form a decision-making state $s = (l, t, I = 0)$
7:         **for** $a \in A$ **do**
8:             Calculate $Q(s, a)$ by Equation (2.8), the optimal values of the dependent future states are read from the hash table $V^*$
9:         **end for**
10:         Derive the optimal policy for decision-making state $s$, $\pi(s)$, by Equation (2.10)
11:         Calculate the optimal value for state $s$, $V^*(s)$, by Equation (2.11)
12:
13:         Form a non-decision-making state $s = (l, t, I = 1)$
14:         Calculate $Q(s, .)$ by Equation (2.9), the optimal values of the dependent future states are read from the hash table $V^*$
15:         Calculate the optimal value for state $s$, $V^*(s)$, by Equation (2.11)
16:     **end for**
17: **end for**
18: return the $V^*$ and $\pi$
___

### 2.3. Extracting parameters from data

In the dataset we used in this research, we have GPS trajectories for both the empty and occupied trips. We now introduce how to extract the parameters we used in the state transition from the dataset.

#### 2.3.1. Order matching probability $P_{order\_match}$

The order matching probability estimates the probability at which a vacant taxi can be matched to a passenger when the taxi is cruising, including staying, or waiting at grid $l_a$. As we have mentioned before, the purposes for introducing waiting and staying in this work are different. In addition to six actions which allow an e-hailing driver to move into one of the six neighboring grids, the action staying gives the driver extra flexibility in choosing to stay and cruise within the current grid due to some potential benefits, such as a relatively high order matching probability in the current grid, a possibly high cost to move into neighboring grids vacantly, etc. Actually, as we have listed in Table (1), there are several studies in the literature that have already included the action staying into the action space, such as (Rong et al., 2016), (Verma et al., 2017), and (Lin et al., 2018). Thus, the way to calculate the order matching probability for staying is the same as the way to calculate the order matching probability for cruising into one of the six neighboring grids. In other words, the order matching probability for a driver just entering the grid from one of the six neighboring grids is supposed to be the same as the order matching probability for a driver who was in the grid and chose to stay in the grid.

Waiting, different from staying and other six actions which allow the driver to move into one of the six neighboring grids, is included into the action space based on the observation that sometimes a driver will choose to stop cruising and simply to wait statically for passenger requests to come in, especially when the driver is around downtown or transportation terminals. The action waiting was previously included in the action space in (Gao et al., 2018).

We thus approximate the order matching probabilities for cruising and waiting separately. We say a driver is waiting for a passenger request whenever the driver's traveling distance is less than 200 meters for a 3-minute interval. To rule out some unrealistic waiting actions, such as a driver being stuck in traffic, we further limit the possible locations for waiting to be the places around subway stations, bus terminals, airports, and some famous tourism attractions. For cruising, the order matching probability can be approximated as the ratio of the number of times that a taxi is matched to a passenger in grid $l_a$ while cruising, denoted as $n_{order\_match\_cruising}(l_a)$, to the number of times that the grid $l_a$ is passed by an empty taxi while cruising, denoted as $n_{passby\_cruising}(l_a)$. For waiting, the order matching probability can be approximated as the ratio of the number of times that a taxi is matched to a passenger in grid $l_a$ while waiting, denoted as $n_{order\_match\_waiting}(l_a)$, to the number of times that empty taxis have waited in the grid $l_a$, denoted as $n_{passby\_waiting}(l_a)$.

$$P_{order\_match}(l_a) = \begin{cases} \dfrac{n_{order\_match\_cruising}(l_a)}{n_{passby\_cruising}(l_a)} & cruising \\ \dfrac{n_{order\_match\_waiting}(l_a)}{n_{passby\_waiting}(l_a)} & waiting \end{cases} \tag{2.12}$$

### 2.3.2. Pick-up probability $P_{pickup}$

The pick-up probability $P_{pickup}(l_a, l'')$ measures the likelihood of picking up a passenger at grid $l''$ when the request sent from the passenger was matched to the driver at grid $l_a$. This parameter can be estimated as the ratio of the the number of passenger pick-ups in grid $l''$ which were matched to drivers in grid $l_a$, denoted as $n_{pickup}(l_a, l'')$, to $n_{order\_match}(l_a)$, which is the summation of $n_{order\_match\_cruising}(l_a)$ and $n_{order\_match\_waiting}(l_a)$.

$$P_{pickup}(l_a, l'') = \frac{n_{pickup}(l_a, l'')}{n_{order\_match}(l_a)} \tag{2.13}$$

### 2.3.3. Destination probability $P_{dest}$

The destination probability $P_{dest}(l'', l''')$ measures the likelihood of the destination of the passenger being grid $l'''$ when the passenger was picked up in grid $l''$. This parameter can be estimated by dividing the number of trips ending in grid $l'''$ which originated from grid $l''$, denote as $n_{dest}(l'', l''')$, by the total number of pick-ups in grid $l''$, denoted as $n_{pickup}(l'')$.

$$P_{dest}(l'', l''') = \frac{n_{dest}(l'', l''')}{n_{pickup}(l'')} \tag{2.14}$$

### 2.3.4. Order matching probability while on trip $P_{match}$

As we have mentioned before, there is a probability at which the driver will receive a request when she is on the trip to transport the current passenger to the destination. We denote this order matching probability while on trip as $P_{match}$. This probability can be estimated by dividing the number of occupied trips among which there is at least one request received by the driver before the driver reaching the destination $l'''$ while the origin is $l''$, denoted as $n_{match}(l'', l''')$, by the total number of occupied trips ending in grid $l'''$ and originating in grid $l''$, denoted as $n_{trips}(l'', l''')$.

$$P_{match}(l'', l''') = \frac{n_{match}(l'', l''')}{n_{trips}(l'', l''')} \tag{2.15}$$

### 2.3.5. Driving time $t_{drive}$ and driving distance $d_{drive}$

The driving time $t_{drive}(l, k)$ and the driving distance $d_{drive}(l, k)$ denote the estimated driving time and driving distance from grid $l$ to grid $k$, respectively. Here we simply take the average of all driving times from grid $l$ to grid $k$ as an approximation of the $t_{drive}(l, k)$. Similarly, the driving distance is calculated by taking the average of all driving distances between grid $l$ and grid $k$.

### 2.3.6. Taxi fare $f$

The taxi fare $f(l'', l''')$ denotes the estimated gross revenue that a driver can earn by transporting a passenger from grid $l''$ to her destination grid $l'''$. Here we take the average of all the fares of the occupied trips which are from grid $l''$ to grid $l'''$ as a proxy of the real taxi fare from grid $l''$ to $l'''$.

### 2.3.7. Seeking time $t_{seek}$ and seeking distance $d_{seek}$

The seeking time $t_{seek}(l_a)$ and the seeking distance $d_{seek}(l_a)$ denote the estimated seeking time and seeking distance within grid $l_a$, respectively. From the field data, the distribution of the seeking time in each grid was extracted and is shown in Figure (3). The median of the distribution of the seeking time is approximately 45 seconds. Since the time step size is 1 minute in this work, thus we simply take the seeking time as 1 minute. Considering the average speed of seeking trips (around 300 meters/minute), the seeking distance is taken as 300 meters for each grid. Note that the seeking distance is zero when the driver chooses to wait.
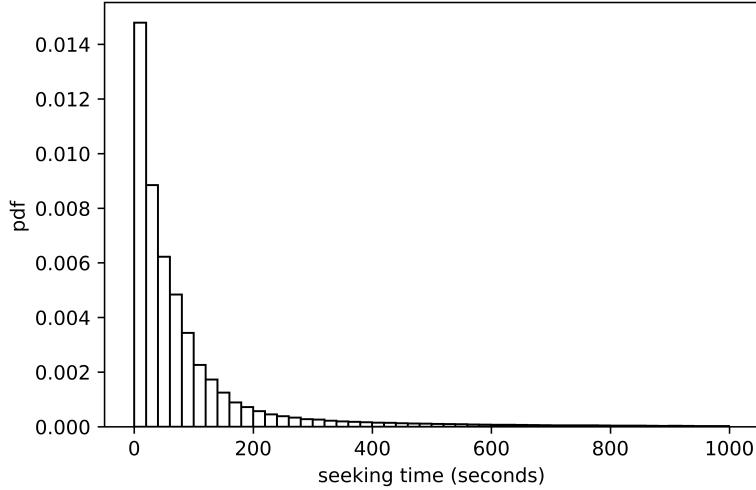
Figure 3: Distribution of seeking time

### 2.3.8. The fuel consumption coefficient $\alpha$

$\alpha$ estimates the fuel consumption and other operating cost per unit distance during driving. In the literature, the value of $\alpha$ is typically assumed to be a known constant based on some common knowledge (Rong et al., 2016; Verma et al., 2017; Lin et al., 2018; Yu et al., 2019). This assumption, however, can easily result in a gap between the reward function in MDP models and the reward function of real drivers, meaning that drivers will not follow the optimal policy since the reward function used to derive the optimal policy is not the reward function of real drivers. To more appropriately determine $\alpha$, we opt for the inverse reinforcement learning (IRL) approach. IRL is a powerful technique to disclose the underlying reward function based on the observed behaviors, especially in the field of imitation learning where an agent's behavior is observed by a learner who tries to imitate the agent (Ziebart et al., 2008).

A linear programming formulation of IRL in finite state spaces was first proposed in (Ng and Russell, 2000), where an extension to large state spaces was also made possible by adopting a linear function approximation. In the context of understanding the observed behaviors, Liu et al. (2013) argued that a mixed integer linear programming formulation can be more revealing. Noticing the similarity between the two approaches, we take the mixed integer linear programming formulation for simplicity. To make the paper self-explanatory, we provide a brief introduction of mixed integer linear programming formulation. Interested readers are referred to (Liu et al., 2013) for detailed explanations.

In the linear programming formulation, the underlying reward function $R(s, s')$ from a state $s$ to another state $s'$ is expressed as a linear combination of some simple known reward functions $\phi_i$s, i.e.,

$$R(s, s') = \alpha_1 \phi_1(s, s') + \alpha_2 \phi_2(s, s') + \cdots + \alpha_n \phi_n(s, s') \tag{2.16}$$

As an example, $\phi_i(s, s')$ can be either the fare a driver can collect from $s$ to $s'$ or the distance a driver traveled from $s$ to $s'$. The former is considered to be positive while the latter is negative. For each simple reward function $\phi_i$, the optimal value function $V^{\phi_i}$ can be derived by solving the MDP with $\phi_i$ as the reward function. Due to linearity, the optimal value function under the underlying reward function can be calculated as

$$V^R = \alpha_1 V^{\phi_1} + \alpha_2 V^{\phi_2} + \cdots + \alpha_n V^{\phi_n} \tag{2.17}$$

The optimal policy $\pi^R$ is also derived when solving the MDP. The objective of the mixed integer linear programming IRL is to minimize the difference between the optimal policy $\pi^R$ and the observed policy $\pi^O$, i.e., minimize$\sum_{s \in S}[\pi^R(s) \neq \pi^O(s)]$. After some mathematical manipulations, a mixed integer linear

programming formulation is formed as

$$\text{minimize} \quad \sum_{s \in S} C_s \tag{2.18}$$

$$\text{s.t.}$$

$$\alpha_i \geq 0 \tag{2.19}$$

$$\sum_{s'} P_{\pi^O}(s, s')(R(s, s') + \gamma V^{\pi^R}(s'))$$

$$- \sum_{s'} P_a(s, s')(R(s, s') + \gamma V^{\pi^R}(s')) + M \times C_s \geq 0 \tag{2.20}$$

where $C_s$ is a binary variable and takes value of 0 or 1 and $M$ is an arbitrarily large number. $\sum_{s \in S} C_s$ denotes the number of states where the observed policy and the optimal policy differ. The objective is thus to minimize the difference between $\pi^R$ and $\pi^O$. Constraint (2.19) restricts the weighting parameter to be nonnegative. In the last constraint (2.20), the first term $\sum_{s'} P_{\pi^O}(s, s')(R(s, s') + \gamma V^{\pi^R}(s'))$ is the optimal value at state $s$ following the observed policy, and the second term $\sum_{s'} P_a(s, s')(R(s, s') + \gamma V^{\pi^R}(s'))$ is the value following other policies. Thus, we expect $\sum_{s'} P_{\pi^O}(s, s')(R(s, s') + \gamma V^{\pi^R}(s')) \geq \sum_{s'} P_a(s, s')(R(s, s') + \gamma V^{\pi^R}(s'))$ to hold. In reality, however, this constraint can be violated, and thus we add a large positive number to the violated constraint to keep (2.20) hold. $P_\pi(s, s')$ is the probability of the transition $s \rightarrow s'$ following the policy $\pi$. In this work, $P_\pi(s, s')$ or $P_a(s, s')$ has been demonstrated in Figure (2).

*Remark.* To identify an overall reward function (i.e., determining $\alpha_i$s), we assume part of the observed policy is optimal or near optimal. Different from simply assuming a reward function based on common knowledge in the literature (Rong et al., 2016; Verma et al., 2017; Lin et al., 2018; Yu et al., 2019), we argue that the assumption used here is relaxed to some degree. Real drivers, at least part of them, are deemed to be intelligent and experienced and exhibit optimal or near optimal strategies. For a state visited by multiple drivers for multiple times, we believe the most frequently taken action carries useful information about the optimal strategy in this state and reflects the crowd wisdom. After uncovering the underlying reward function, the purpose of solving the MDP and thus deriving the optimal policy is to complete and correct the optimal policy. The incompleteness and noise in the observed policy stem from the following aspects: (1) For a real-world problem with a huge number of states (in our case study in Section 4 the number of decision-making states is $6,421 \times 180 = 1,115,780$), a considerable portion of the states are not visited or at least not frequently visited; (2) Even in states with sufficient data (i.e., enough actions chosen by agents in this state), the most frequently chosen action which is taken as the observed policy in this state can still differ from the subsequent derived optimal policy, due to behavioral inconsistency; and (3) the observed policy, which is assumed to be deterministic, can be ambiguous when several actions share similar frequency in some states.

## 2.4. Numerical example

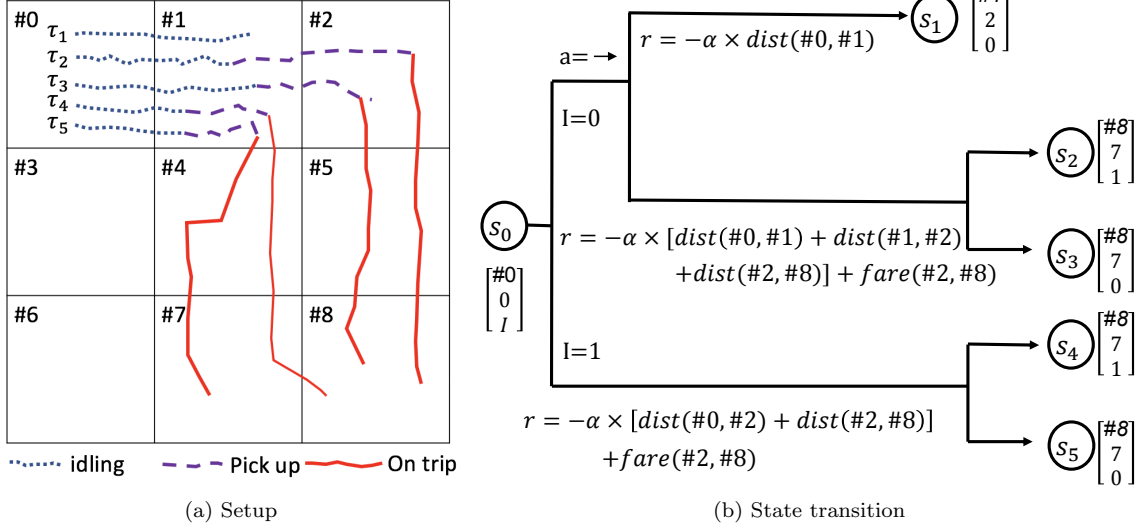

(a) Setup            (b) State transition

Figure 4: Setup of the small grid world example and the corresponding state transition

To illustrate the Markov Decision Process of e-hailing drivers, we use a 3 by 3 grid world numerical example, as shown in Figure (4a).

Suppose now we have the following five trajectories.

1. $\tau_1 = (\#0, 0, 0) \xrightarrow{\text{idling}} (\#1, 2, 0)$

2. $\tau_2 = (\#0, 0, 0) \xrightarrow{\text{idling}} (\#1, 2, 0) \xrightarrow{\text{pickup}} (\#2, 3, 0) \xrightarrow[\text{no match}]{\text{ontrip}} (\#8, 7, 0)$

3. $\tau_3 = (\#0, 0, 0) \xrightarrow{\text{idling}} (\#1, 2, 0) \xrightarrow{\text{pickup}} (\#2, 3, 0) \xrightarrow[\text{match}]{\text{ontrip}} (\#8, 7, 1)$

4. $\tau_4 = (\#0, 0, 0) \xrightarrow{\text{idling}} (\#1, 2, 0) \xrightarrow{\text{pickup}} (\#1, 3, 0) \xrightarrow[\text{no match}]{\text{ontrip}} (\#8, 7, 0)$

5. $\tau_5 = (\#0, 0, 0) \xrightarrow{\text{idling}} (\#1, 2, 0) \xrightarrow{\text{pickup}} (\#1, 3, 0) \xrightarrow[\text{no match}]{\text{ontrip}} (\#7, 6, 0)$

Each element in the trajectory is a tuple consisting of three items, namely, the grid index, current time, and a status indicator showing if the driver has been matched to another order during the trip. For example, $(\#0, 0, 0)$ basically states that the driver is at grid $\#0$ at time 0, and the driver has not been matched to any order before she finished the previous trip.

All five trajectories started in grid $\#0$, and then the driver moved into grid $\#1$ during idling. After the driver searched the grid $\#1$, there are two possible outcomes: either the driver finds an order match or the driver fails to find any e-hailing order. If the driver fails to get an order match after searching, the driver will move into other grids to find another order or the driver will stop working. To simplify the demonstration, we simply assume the trajectory $\tau_1$ ends in grid $\#1$. For other four trajectories, the driver managed to find an e-hailing order in grid $\#1$. Based on this piece information, we can calculate the probability of finding an e-hailing order in grid $\#1$ as $P_{order\_match}(\#1) = \frac{n_{order\_match}(\#1)}{n_{passby}(\#1)} = \frac{4}{5} = 80\%$.

After confirming an order match in grid $\#1$, the driver drives into grid $\#2$ to pick up the passenger in trajectories $\tau_2$ and $\tau_3$ and stays within grid $\#1$ to pick up the passenger in trajectories $\tau_4$ and $\tau_5$, respectively. Thus, the pick-up probability can be calculated as $P_{pickup}(\#1, \#1) = \frac{n_{pickup}(\#1, \#1)}{n_{order\_match}(\#1)} = \frac{2}{4} = 50\%$ and $P_{pickup}(\#1, \#2) = \frac{n_{pickup}(\#1, \#2)}{n_{order\_match}(\#1)} = \frac{2}{4} = 50\%$.

14

When the driver picks up the passenger in grid #1, as illustrated in trajectories $\tau_4$ and $\tau_5$, the passenger's destination is grid #8 in $\tau_4$ and grid #7 in $\tau_5$, respectively. Thus, the destination probability can be calculated as $P_{dest}(\#1, \#7) = \dfrac{n_{dest}(\#1, \#7)}{n_{pickup}(\#1)} = \dfrac{1}{2} = 50\%$ and $P_{dest}(\#1, \#8) = \dfrac{n_{dest}(\#1, \#8)}{n_{pickup}(\#1)} = \dfrac{1}{2} = 50\%$.

In trajectories $\tau_2$ and $\tau_3$, the driver drives to grid #2 to pick up the passenger, and the passenger goes to grid #8. During the trip, the driver has a $P_{match}$ of receiving a new order before she arrives at the destination of the passenger. The order matching probability while on trip can thus be calculated as $P_{match}(\#8) = \dfrac{n_{match}(\#8)}{n_{trips}(\#8)} = \dfrac{1}{2} = 50\%$.

Based on the probabilities calculated above, an example of state transition is presented in Figure (4b). To make the state transition consistent with the five trajectories, we suppose the driver is initially in grid #1 with a status indicator $I$, i.e., the driver is in state $s_0 = (\#0, 0, I)$. If $I = 0$, meaning that the driver needs to seek for an e-hailing order, the driver drives into grid #1 and seeks for e-hailing orders in the grid. There are two possible outcomes associated with this case.

1. The driver fails to find any e-hailing order in grid #1. The driver will end up in state $s_1 = (\#1, 2, 0)$ and receive a negative reward $-\alpha \times dist(\#0, \#1)$, which is actually the fuel cost. This outcome happens with probability $p_1 = 1 - P_{order\_match}(\#1) = 1 - 80\% = 20\%$

2. The driver successfully finds an e-hailing order in grid #1. For the purpose of demonstration, we assume the driver goes to grid #2 to pick up the passenger, and the destination of the passenger is grid #8. The probability of this outcome is $p_2 = P_{order\_match}(\#1) \times P_{pickup}(\#1, \#2) \times P_{dest}(\#2, \#8) = 80\% \times 50\% \times 100\% = 40\%$. The driver will receive a total reward of $r = fare(\#2, \#8) - \alpha \times [dist(\#0, \#1) + dist(\#1, \#2) + dist(\#2, \#8)]$ by completing this ride. During the trip, the driver may have a probability $P_{match}(\#8) = 50\%$ of getting a new request before she arrives at the destination of the previous passenger. Hence, there are two possible subbranches from this outcome.

   (a) If the driver is matched to a new request before she drops off the previous passenger, then the driver will end up in state $s_2 = (\#8, 7, 1)$. The probability of the occurrence of this subbranch is $p_2 \times P_{match}(\#8) = 40\% \times 50\% = 20\%$.

   (b) if the driver fails to be matched to another request while on trip, the driver will then end up in $s_3 = (\#8, 7, 0)$. This subbranch occurs with a probability $p_2 \times (1 - P_{match}(\#8)) = 40\% \times (1 - 50\%) = 20\%$

For the sake of the completeness of the state transition, the other two subbranches associated with $I = 1$ are also displayed in Figure (4b). These two subbranches are quite self-explanatory, and thus the detailed discussion will be omitted.

## 3. Sequential MDPs for multiple agents

Note that the deterministic policy derived is only applicable when there is one agent following the policy. Otherwise there can be local competition among e-hailing drivers since several drivers may be guided into the same grid. We thus need to address the competition among e-hailing drivers if there are multiple idling e-hailing drivers being present in the same region within a short time interval. Lin et al. (2018) proposed a contextual multi-agent reinforcement learning approach in which the multi-agent effect is captured by attenuating the reward through an averaging fashion. Zhou et al. (2018) employed a simple discounting factor $\dfrac{1}{n}$ to update the order matching probability when the $(n+1)^{th}$ taxi is being guided to a road if there are already $n$ taxis going to that road. The discounting factor proposed $\dfrac{1}{n}$ is effective in the sense that it makes the order matching probability smaller for subsequent taxis following the policy. However, the simple discounting factor may underestimate the order matching probability since the effect of the number of orders in each grid was neglected. In other words, except the effect of the number of drivers being guided into a grid, there is an underlying correlation between the decrease in the order matching probability and the number of orders in that grid. Here we use an example to illustrate the existence of the aforementioned correlation. We suppose an e-hailing driver is guided into

grid $l$ with an order matching probability 50%. We consider two extreme scenarios: (1) there was 1 order emerging in grid $l$ and (2) there were infinite orders emerging in grid $l$. After one driver is guided into grid $l$, for a second driver, the order matching probability in grid $l$ is supposed to decrease substantially in the first scenario while almost keeps the same in the second scenario. The rationale underlying this argument is that compared to a grid with a smaller number of orders, a grid with a larger number of orders is capable of accepting more cruising drivers while still maintain a relatively acceptable level of order matching probability.

To incorporate this correlation, we develop a dynamic adjustment strategy. Before formally providing the form of the strategy, we list four intuitive observations: (1) The order matching probability for the first driver being guided into grid $l$ is simply $P_{order\_match}(l)$; (2) The order matching probability for the $n^{th}$ driver being guided into grid $l$ decreases with $n$, meaning that the order matching probability is getting smaller when there are more drivers cruising vacantly in grid $l$; (3) For the $n^{th}$ driver, the order matching probability increases with the number of orders in grid $l$, meaning that a grid with a larger number of orders is able to accept more cruising drivers; (4) Under the extreme scenario where there are infinitely many orders in grid $l$, the order matching probability keeps its level at $P_{order\_match}(l)$ regardless of the number of drivers being guided into grid $l$, as long as it is finite. Based on these four observations, we postulate that the order matching probability of the $n^{th}$ driver in grid $l$ takes the exponential form, i.e.,

$$Pr(l,n) = P_{order\_match}(l) \times e^{-\frac{\beta}{\#orders(l)} \times (n-1)} \tag{3.1}$$

where $\#orders(l)$ is the number of orders in grid $l$ and $\beta$ is a parameter to be determined.

To calibrate the strategy for the $n^{th}$ driver, the order matching probability $Pr(l,n)$ is required. Note that $Pr(l,1) = P_{order\_match}(l)$, then Equation (3.1) can be written as

$$\frac{Pr(l,n)}{Pr(l,1)} = e^{-\frac{\beta}{\#orders(l)} \times (n-1)} \tag{3.2}$$

To utilize linear regression, we further take the logarithm and then the reciprocal of both sides of Equation (3.2), and thus Equation (3.2) can be rewritten as

$$\frac{1}{log\frac{Pr(l,n)}{Pr(l,1)}} = -\frac{1}{\beta \times (n-1)} \times \#orders(l) \tag{3.3}$$

Denoting the left hand side of Equation (3.3) as $prob_n(l)$, the purpose of the calibration is simply to verify the existence of the linear correlation between the two variables $\#orders$ and $prob_n$ and determine the parameter $\beta$. $\#orders(l)$ is simply the number of historical orders in grid $l$. The order matching probability for the $n^{th}$ driver entering grid $l$ can be calculated as follows. We use 18 10-minute intervals to split the 3-hour morning peak used in this study. For each interval, the number of orders within the interval is counted, and it is a success if the counted number is not less than $n$. The probability of success across 18 intervals is $Pr(l,n)$. Calculating $prob_n(l)$ and $\#orders(l)$ for all grids, samples of $prob_n(l)$ and $\#orders(l)$ are obtained.

Ideally, the calibration can be run for any integer value of $n$. However, due to the relatively small size of the grid and the finite number of idling drivers, $n$ cannot be taken as a very large number. To get an appropriate upper bound of $n$ in the calibration, we simply count the number of cases in which there are $n$ drivers idling in a grid within a time interval. The distribution of the number of cases versus the number of idling drivers in a grid is presented in Figure (5). For example, there are more than 25,000 cases in which there are only one idling driver in a grid within a 10-minute time interval. The trend is decreasing, indicating that the number of cases for more drivers idling in a grid within a time interval is less, which is as expected. To obtain statistically meaningful results in the calibration, we need as many samples as possible. Thus, here we choose to do the calibration for $n$ up to 4.
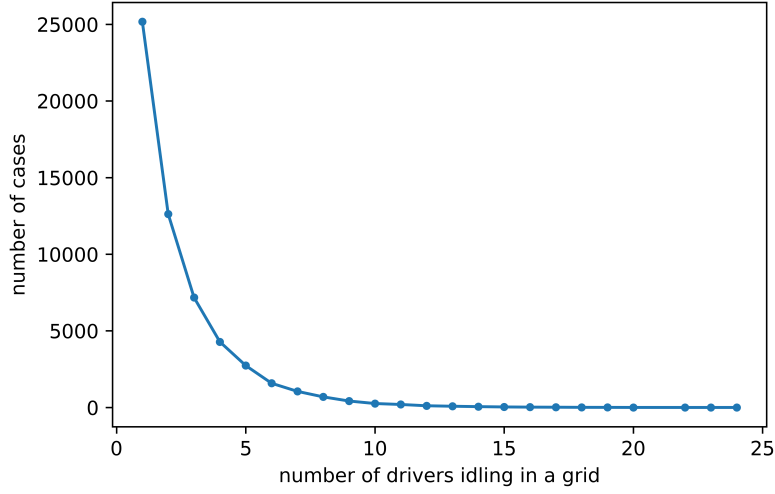
16

Figure 5: Distribution of number of cases

We obtain three lists of $prob_n$ and three lists of $\#orders$ for $n = 2, 3$, and $4$, although the three lists $\#orders$ for different $n$ are identical. Noticing that here we only have one parameter $\beta$, we concatenate the three lists of $prob_n$ simply into a long list $prob = [prob_2, prob_3, prob_4]$ and concatenate the three lists of $\#orders$ as a long list $\#orders\_long = [\#orders, \dfrac{\#orders}{2}, \dfrac{\#orders}{3}]$. Then we run a linear regression through origin of the list $prob$ over $orders\_long$. Here we choose the linear regression through origin because the model is naturally linear without intercept, as shown in Equation (3.3), stemming from the requirement that the order matching probability for the first driver in a grid $l$ is $Pr(l, 1) = P_{order\_match(l)}$. The regression result suggests the slope $-\dfrac{1}{\beta} = -0.0842$ with p-value less than 0.001, indicating that the slope is significant. The R-square value is determined as 65%, indicating the fitted linear model is able to explain 65% of the variability and thus the adoption of a linear regression is reasonable. To visually show the goodness of the fitting, we substitute the determined value of $\beta$ into Equation (3.1) and plot the fitted curve together with the data for $n = 2, 3$, and $4$ in Figure (6).
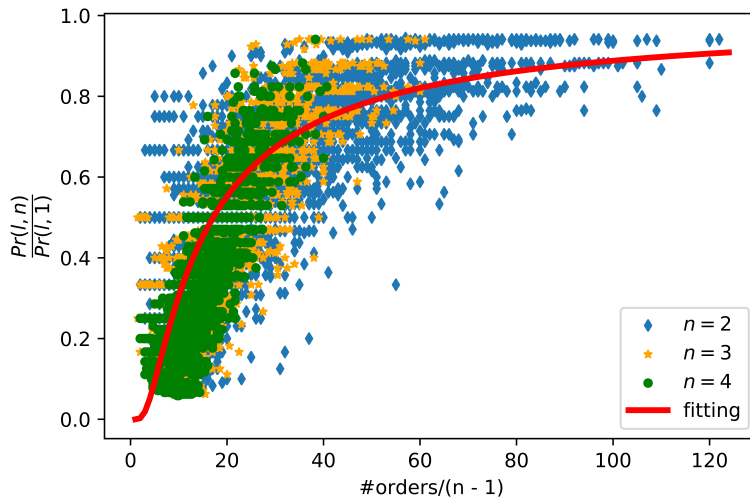


Figure 6: Exponential fitting

Based on the calibrated dynamic adjustment strategy, we can now build sequential MDPs for multiple agents and then derive one optimal policy for each agent sequentially. Recall that the dynamic adjustment

17

strategy actually attenuates the order matching probability for the $n^{th}$ driver to be guided into grid $l$ when there are $n-1$ drivers idling in grid $l$. Thus, the general MDP framework proposed in the previous section can directly be applied to the $n^{th}$ driver with a table of relatively lower order matching probabilities. In other words, there is one MDP model for each agent when multiple agents coexist, and the main distinction among these MDP models is the order matching probability (i.e., order matching probabilities of the $n^{th}$ driver are smaller than that of the $(n-1)^{th}$ driver). The optimal policy for each agent can then be obtained by solving her MDP model.

## 4. Case study

Nowadays, GPS-enabled devices are ubiquitously used in different types of applications. Especially for drivers, GPS devices can not only help them navigate but also record the real-time location and speed of the vehicle. Hence, a large amount of GPS trajectories, representing sequences of time-stamped geographical points, have been collected and can be a valuable asset for people to understand and tackle real-world traffic problems (Di et al., 2010, 2017; Shou and Di, 2018).

In this research we use large-scale real-world historical GPS traces collected in Beijing in the morning peak, i.e., (7 AM, 10 AM), of the first 3 weekdays in November 2017 by Didi Chuxing, China's leading ride-hailing company. The dataset contains recorded GPS traces of 44,160 e-hailing vehicles and 158,784 e-hailing orders. Whenever an e-hailing driver is online, i.e., she has turned on the e-hailing application, one data point is sampled every 3 seconds. Each data point contains the vehicle's location (i.e., longitude and latitude), current timestamp, the fare of the current occupied trip, if applicable, and a status indicator to record the taxi's current operating state, which includes idle, after matching before pick-up, waiting at the pick-up location, and on trip.

Based on the characteristics of the spatial distribution of the passenger pick-up spots and passengers' destinations, we construct a bounding box within the six ring road to cover the city area in which 90% of the e-hailing orders are preserved. The e-hailing orders which fall outside the bounding box will be disregarded. A hexagonal grid world setup is then adopted, and the city area within the bounding box is split into 6,421 hexagonal grids with the length of the diagonal of a hexagon of approximately 700 meters.

Now we use the IRL technique to uncover the general reward function for all drivers, i.e., to derive the parameter $\alpha$. The observed policy $\pi^O(s)$ in each state $s$ is simply the most frequently taken action by all drivers in that state. Applying the aforementioned IRL technique with two known reward functions (i.e., fare $\phi_1(s, s')$ and traveling distance $-\phi_2(s, s')$) and setting $\alpha_1 = 1$ (i.e., assuming the driver will earn all the fare), we obtain $\alpha = \alpha_2 = 0.64$. In other words, the coefficient of fuel consumption and other operating costs per unit distance is 0.64 Chinese Yuan. The general reward function applicable to all drivers is $R(s, s') = \phi_1(s, s') - 0.64 \times \phi_2(s, s')$.

Different from other public transportation modes, including buses and subways, which are operated according to a fixed schedule and a predefined route, e-hailing drivers are free to choose their own actions after completing a ride, resulting in a discrepancy in drivers' income. In general, an experienced driver is familiar with the city where she operates the vehicle and is usually aware of where to go after completing a ride in order to get the next request as soon as possible. For example, an experienced driver knows where and when passenger demands will be high near some attractions, hotels, or transportation terminals. In practice, however, there is no guarantee that all drivers are experienced and have a good judgment of where to go next. In particular, the popularity of e-hailing applications has lowered the entry barrier of becoming a driver and brought a tons of rookie into the e-hailing market. To gain some basic understanding of the performance of e-hailing drivers, we adopt two metrics, namely, the rate of return and the utilization rate.

**Definition 4.1.** (Rate of return) An e-hailing driver's rate of return on one day is defined as the ratio of the driver's net income (i.e., gross income minus the operating cost) to the driver's working time.

**Definition 4.2.** (Utilization rate) The utilization rate of an e-hailing vehicle is defined as the ratio of the time spent on carrying a passenger to the total operating time of the vehicle.
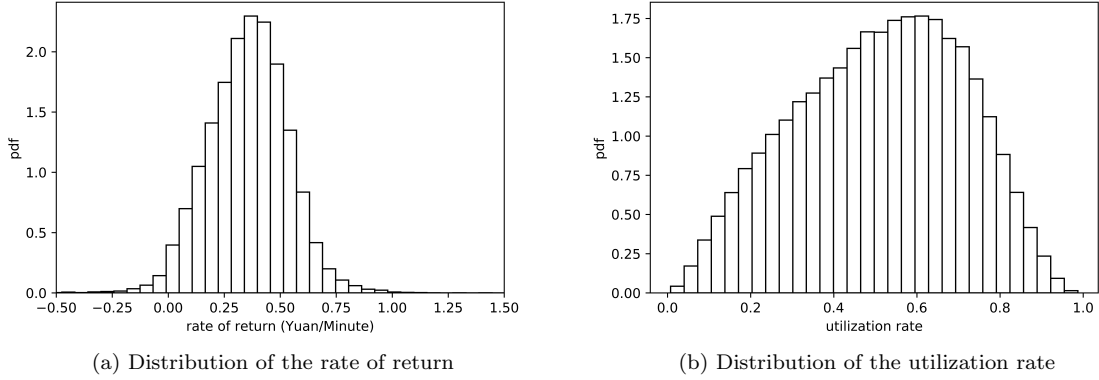
(a) Distribution of the rate of return        (b) Distribution of the utilization rate

Figure 7: Distribution of the rate of return and utilization rate from the field data

The rate of return measures a driver's earning ability per unit time. Thus, compared with the total income, which can be largely influenced by a driver's working time, the rate of return is a better metric to be used to measure the performance of drivers. Figure (7a) presents the probability density function (pdf) of the rate of return of all drivers across morning peaks on the first three weekdays in November 2017. The average rate of return is 0.36 (Yuan/minute). About 80% drivers have a rate of return fall within the range 0.11 to 0.60 (Yuan/minute). Top 10% drivers can reach a rate of return of 0.60 (Yuan/minute) and higher, while the bottom 10% have a rate of return below 0.11 (Yuan/minute). In terms of the utilization rate, real drivers can on average reach 0.51.

All drivers' data was then used to train the MDP model. We will then qualitatively examine the optimal policy derived from our MDP model and conduct numerical experiments to evaluate the effectiveness of the policy.
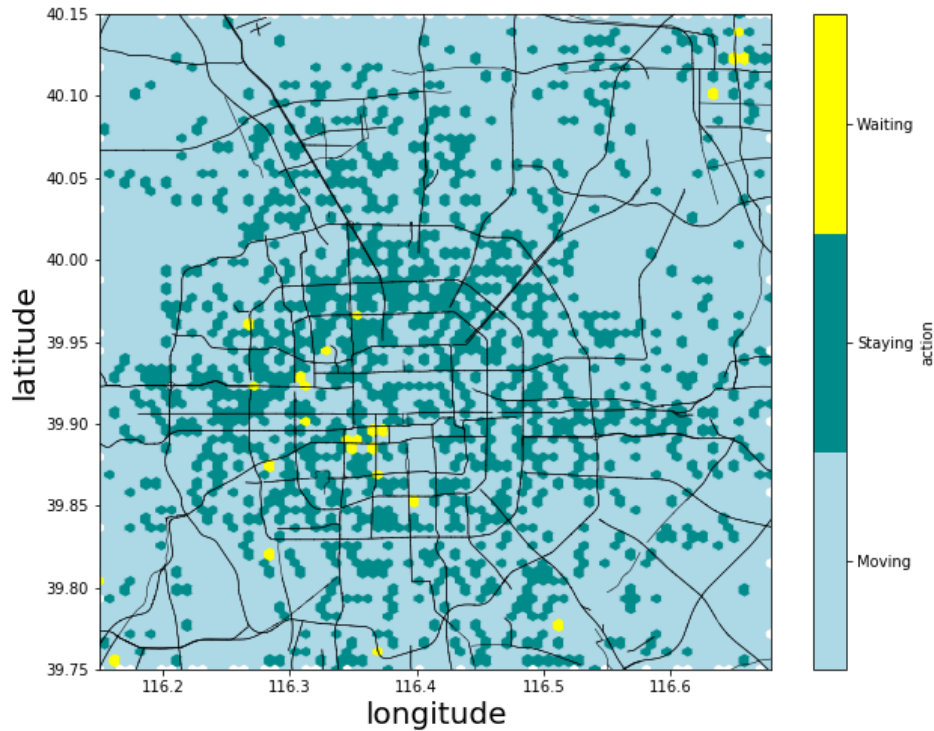
Figure 8: Optimal policy (at the beginning of the morning peak)

Figure (8) presents the optimal policy at the beginning of the morning peak. Light blue color suggests the driver in the current grid to move into one of its neighboring grids to seek for the next potential e-hailing order. The color dark green in a grid stands for staying, meaning that the optimal policy is to stay and cruise within the current grid. The color yellow in a grid means waiting, indicating that the optimal policy for the driver to follow is simply to wait in the current grid. It can be seen that in many grids within the city area (i.e., around the center part of the figure), the optimal policy suggests a driver to stay after completing a ride. In suburban areas, optimal policy usually suggests drivers to move around to some grids with a high probability of receiving a request. Also, there are several places where the probability of receiving a request while waiting is quite high, and the optimal policy advices drivers to wait in these places.
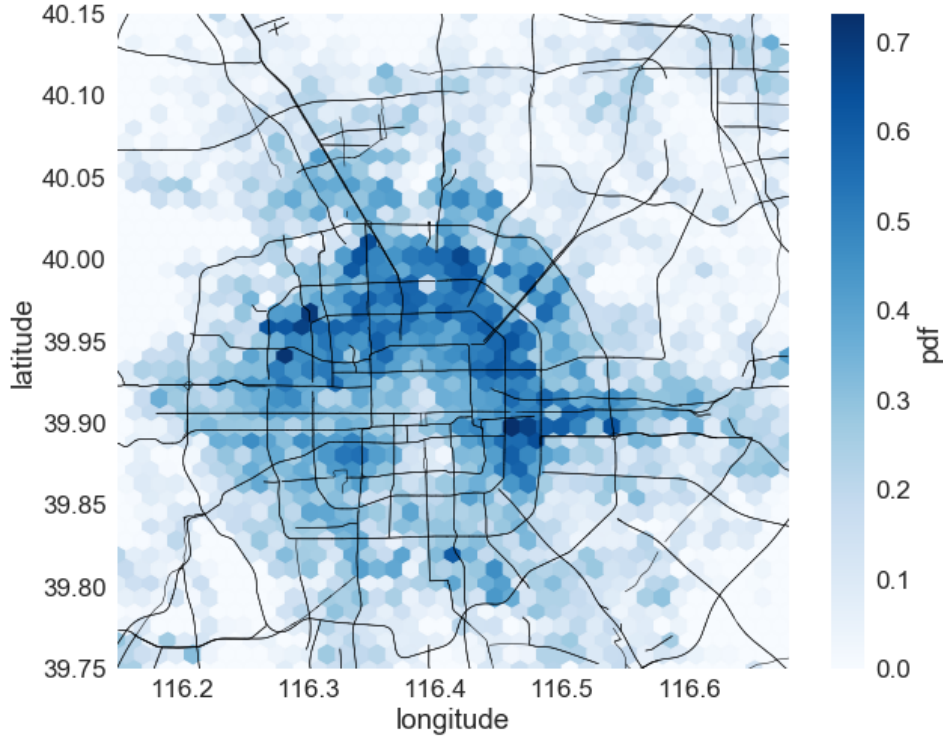
Figure 9: Distribution of the order matching probability $P_{order\_match}$

Figure (9) plots the distribution of the order matching probability. It can be seen that the majority of the color dark blue, indicating a higher order matching probability, is distributed within the fourth ring road, meaning that the order matching probability is higher in the city area. Furthermore, the distribution of the order matching probability agrees well with the distribution of the optimal policy if we compare Figure (9) with Figure (8). The dark green places, indicating staying, in Figure (8) generally overlap with blue or dark blue locations in Figure (9), indicating a relatively high order matching probability. Also, when a driver is in a grid with a low order matching probability (e.g., in the nearly white area in Figure (9)), the driver needs to move around (as shown by the light blue in Figure (8)) to enter a grid with a higher order matching probability. This overlapping essentially indicates that grids with a higher order matching probability is more preferable by the agent, compared with a grid with a lower order matching probability.

### 4.2. Model Evaluation

To evaluate the effectiveness of the optimal policy derived from the MDP model, the performance of an agent under the guidance of the optimal policy is compared with that of an agent following one baseline heuristic, i.e., the local hotspot strategy. The local hotspot strategy essentially suggests the agent to move into grids with a higher demand sequentially and is found to perform the best among three baseline heuristics, namely, random walk, global hotspot, and local hotspot (Yu et al., 2019).

To obtain the performance of an agent according to different policies, a Monte Carlo simulation is conducted. The basic idea of the simulation is to randomly place an agent in one grid at the beginning, and let the agent move around according to the chosen policy. The environment is determined by the parameters extracted in Section (2.3). In particular, every time when the agent is in a grid $l_a$, we sample a probability of finding a request from a binomial distribution with a success probability $p_{order\_match}(l_a)$. We then sample the pick-up grid and drop-off grid from a multinomial distribution determined by the probabilities $p_{pickup}$ and $p_{dest}$, respectively. The driving time and driving distance can be simply obtained from $t_{drive}$ and $d_{drive}$ after the pickup spot and destination have been determined.

The simulation is run for millions of times to obtain robust results. We first adopt two metrics, namely, rate of return and the utilization rate of the vehicle to compare the performance of the agent under the guidance of different policies. We then examine the distribution of the number of completed orders, idling time, service time per order, and the profit per unit time of each order.
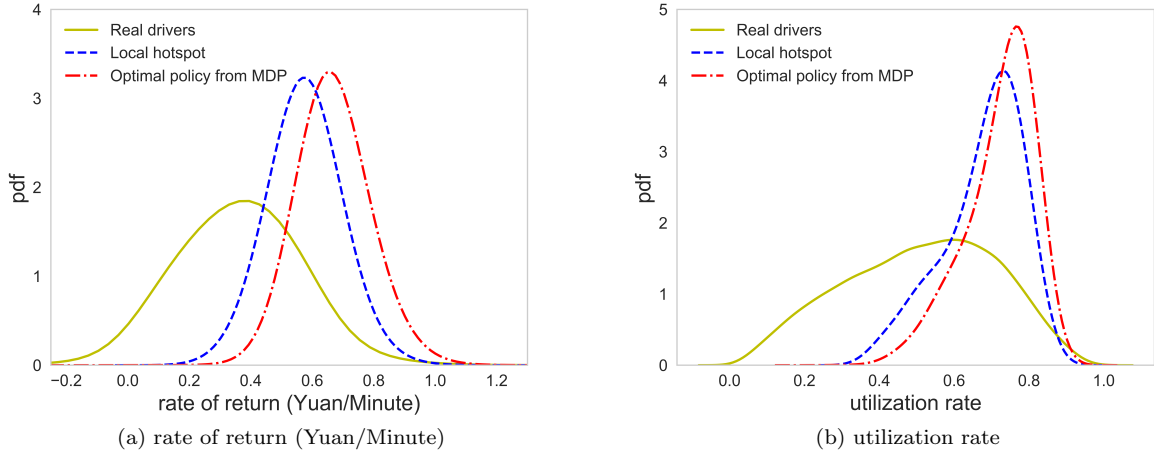


(a) rate of return (Yuan/Minute)
(b) utilization rate

Figure 10: Distribution of the rate of return and the utilization rate of real drivers and the agent

Figures (10) presents the distribution of the rate of return and the utilization rate of the agent following the optimal policy and the local hotspot strategy, respectively. It can be seen that the performance of the agent following the optimal policy is on average better than that of the agent following the local hotspot strategy. In terms of the rate of return, the average value that the agent can reach are 0.67 (Yuan/Minute) and 0.57 (Yuan/Minute) under the guidance of the optimal policy and the local hotspot strategy, respectively, meaning that the optimal policy is able to increase the average rate of return by 17.5% over the local hotspot strategy. The average rate of return of real drivers is 0.36 (Yuan/Minute). In terms of the utilization rate, the average value that the agent can reach are 0.72 and 0.67 by following the optimal policy and the local hotspot strategy, respectively, indicating a 7.5% improvement of the optimal policy over the local hotspot strategy. The average utilization rate of real drivers is around 0.51.

Table 3: Statistics of the comparison between the performance of the agent under different policies

|  | real drivers | agent following the local hotspot strategy | agent following the optimal policy |
|---|---|---|---|
| average rate of return (Yuan/Minute) | 0.36 | 0.57 | 0.67 |
| average utilization rate | 0.51 | 0.67 | 0.72 |
| average number of orders | 6.08 | 8.21 | 8.92 |
| average idling time (Minute) | 47.98 | 27.12 | 22.61 |
| average profit per unit time of each order (Yuan/Minute) | 1.35 | 1.51 | 1.68 |
| average service time per order (Minute) | 16.97 | 14.83 | 14.56 |

(a) Distribution of number of orders        (b) Distribution of idling time

(c) Distribution of profit per unit time of each order        (d) Distribution of service time per order
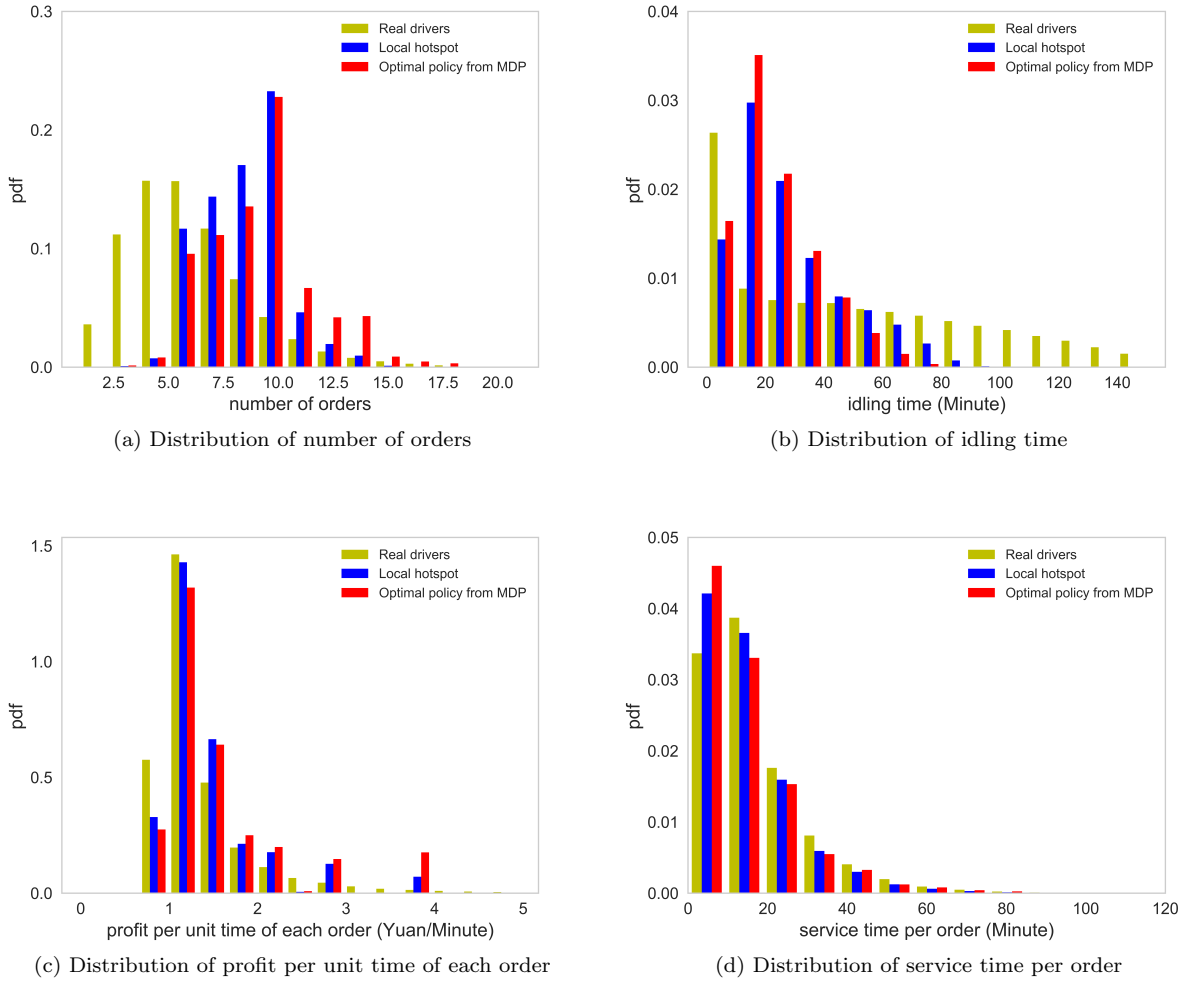
Figure 11: Distributions of the number of completed orders, number of idling cells, the service time per order, and the profit per unit time of orders

Figure (11) presents distributions of the number of completed orders, idling time, service time per order, and the profit per unit time of each order. The corresponding average value of each case is listed in Table (3). There are several observations worth mentioning: (1) On average the agent following the optimal policy can achieve a higher number of completed orders than the agent following the local hotspot strategy. (2) The agent following the optimal policy and the local hotspot strategy can achieve a significantly lower idling time, compared with that of real drivers. Thus, both the optimal policy and the local hotspot strategy are able to help agents find requests faster. (3) The distribution of the profit per unit time essentially says that under the optimal policy, the agent is able to find better orders. Here we say an e-hailing order is better when the profit per unit time of the order is higher. (4) The average service time of orders taken by the agent is shorter than that of real drivers.

The aforementioned observation (4), however, does not necessarily indicate that on average a shorter order is more preferable compared with a longer order. The reasons are as follows. As previously stated, for an agent, grids with relatively higher order matching probability is more preferable, compared with grids with a moderate or low order matching probability. Figure (12) presents the distribution of service time of orders starting in all grids and in grids with a relatively high order matching probability. The average service time of orders starting in all grids and in grids with a relatively high order matching probability are 17 minutes and 18.34 minutes, respectively. In other words, the passenger requests in grids with a relative higher order matching probability on average has a slightly greater service time. Hence, the average service time per order of the agent is supposed to be slightly higher than that of real drivers, which is not consistent with the observation (4).

To explain the inconsistency, we split the 3-hour time interval, i.e., the morning peak, into 6 subintervals of even length, namely, $[0, 30)$, $[30, 60)$, $[60, 90)$, $[90, 120)$, $[120, 150)$, and $[150, 180]$, and then extract the average service time of orders starting in each subinterval, which is shown in Figure (13). The average service time of orders starting in each subinterval decreases as time elapses. In the first subinterval, i.e., around the beginning of the morning peak, the agent actually takes orders with an average service time of 16.88 minutes, which is very close to 16.97 minutes, i.e., the average service time of all orders. This is as expected because the initial position of the agent is randomly chosen, meaning that the distribution of the service time of orders taken by the agent at the beginning of the time interval is supposed to be similar to that of all orders. As time elapses, the average service time of orders taken by the agent decreases because in the simulation we restrict the ending time of an order that an agent can take to be within the 3-hour time interval, and thus the agent to some degree prefers shorter orders, especially when the agent is in the last two subintervals, i.e., around the end of the 3-hour time interval. Low values of the service time of orders in the last two subintervals, i.e., 14.12 minutes and 9.01 minutes largely drag down the overall average service time of orders taken by the agent and caused the aforementioned inconsistency. Here we emphasize that in the simulation, the agent is set to complete the 3-hour time interval, and the restriction imposed by the 3-hour time interval implicitly forces the agent to take relatively shorter orders. All references listed in Table (1) except (Gao et al., 2018) used finite horizon time intervals (Rong et al., 2016; Verma et al., 2017; Lin et al., 2018) or finite pickup-dropoff cycles (Yu et al., 2019) in MDP modeling. Actually, in reality, an e-hailing driver is free to stop working at any time as long as she has achieved her preset goal, such as a nominal income, a personalized utility function, etc. Thus, an MDP with optimal stopping time modeling is a more realistic and efficient model and is left for future research.
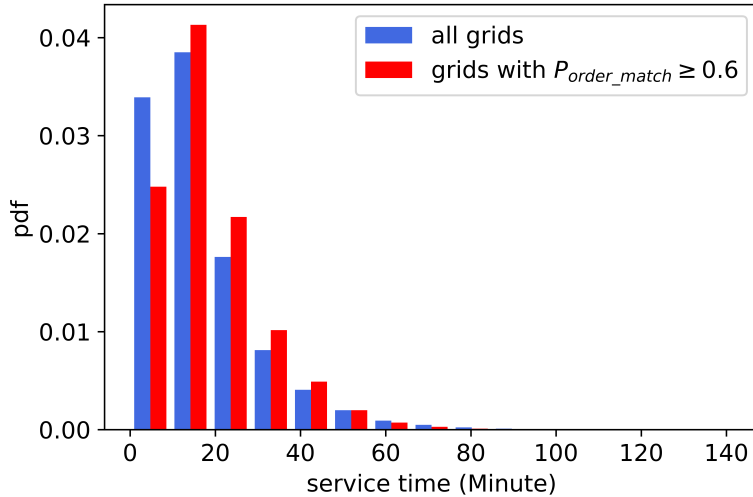


Figure 12: Distribution of service time of orders starting in all grids and in grids with a relatively high order matching probability
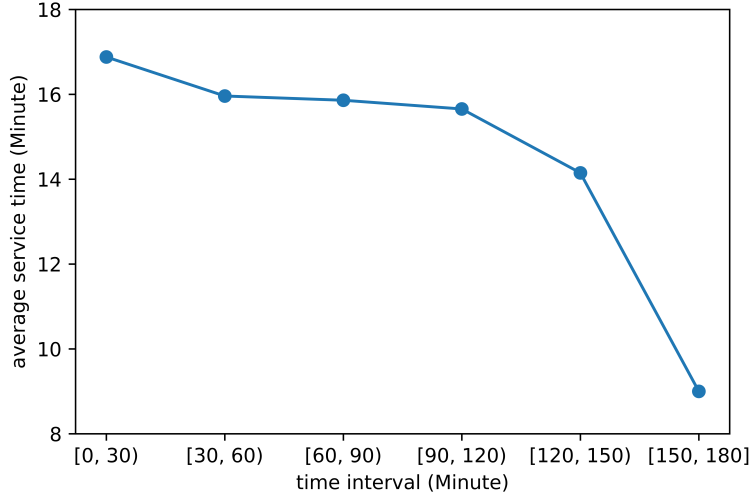
Figure 13: Average service time of orders starting in each subinterval

## 4.3. Supply-demand ratio



(a) Around Xibeiwang (outside the 5th ring road)
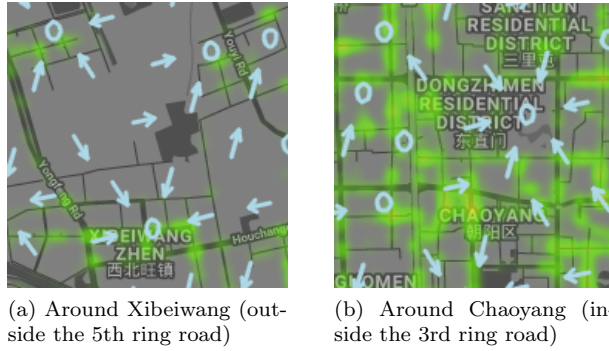


(b) Around Chaoyang (inside the 3rd ring road)

Figure 14: Two zoom-in views of the optimal policy and the distribution of the demand

Figure (14a) presents a zoom-in view of the optimal policy (shown in arrows and circles where arrows suggest the driver to move along the direction denoted by the arrow and the circle denotes staying or waiting) and the distribution of the demand (shown in heatmap) near Xibeiwang, an area outside the 5th ring road. We can see that our optimal policy suggests drivers to stay around the grid where the demand is high to take advantage of the locally high demand. Outside the 5th ring road, the overall demand is not very high and the number of idle drivers, indicating the supply, is also not large, resulting in a locally high order matching probability in high demand areas. Thus, a driver can simply take advantage of the high demand to make more profit.

Figure (14b) presents a zoom-in view of the optimal policy around Chaoyang district, which is located within the 3rd ring road. Different from the consistency between the policy and the distribution of the demand observed in Figure (14a), now it seems the derived optimal policy and the distribution of the demand is not very consistent. In other words, the optimal policy suggests drivers to move away from the areas with a high demand. The main reason is that the optimal policy is supposed to be consistent with the order matching probability, which captures the supply-demand ratio under the assumption that the number of unmatched order is negligible, instead of the real distribution of the demand. Although the order matching probability is calculated from the distribution of the demand and the number of pass-bys of idle drivers, there may exist some shift or even inconsistency between the order matching probability and the distribution of the demand. For instance, a cell with a very high demand, e.g. 100 order matches, may have a 1,000 pass-bys by idle drivers, resulting in a relatively low order matching probability in the cell, which is 0.1 in this example. Hence, it is not surprising that the derived optimal

policy suggests the driver to move away from the cell with a high demand and a low order matching probability. A grid with a high demand may also have a high supply, resulting in a low order matching probability which is not preferable to the agent.

### 4.4. Different optimal policies for multiple agents



(a) Around Dongxiaoying (outside the 5th ring road)

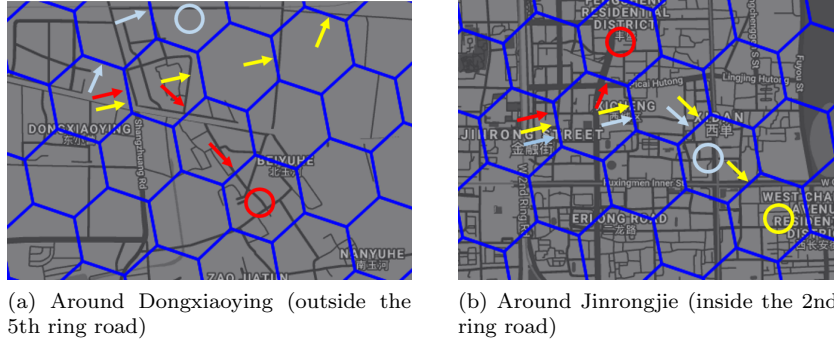(b) Around Jinrongjie (inside the 2nd ring road)

Figure 15: Two zoom-in views of the optimal policy for three agents (red, yellow, and light blue denote the first, second, and third agent, respectively)

Now, we verify the effectiveness of the proposed dynamic adjustment strategy of the order matching probability by recommending the next several optimal actions to take for three agents who are currently in the same grid. We randomly selected two places, namely, one place around Dongxiaoying (outside the 5th ring road, thus in suburban area) and one place around Jinrongjie (inside the 2nd ring road, thus in city area).

Figure (15) presents the recommended optimal actions to take for three agents. In both cases, three agents are being recommended at the same time in the same grid. The results shows that our proposed dynamic adjustment strategy is able to provide different recommended actions for different agents. When three agents are in the same grid around Dongxiaoying, which is located outside the 5th ring road and has a relatively lower number of orders, the strategy guides the first two agents into the same grid but refers the third agent into a different grid because after guiding two agents into the same grid, the order matching probability in that grid drops quickly. The strategy also guides the first two agents into two different directions afterwards due to the competition. When three agents are in the same grid around Jinrongjie, which is located inside the 2nd ring road and has a relatively higher number of orders, the strategy guides all three agents into the same grid. Although the strategy refers the first agent into a different direction soon, the strategy guides the following two agents in almost the same direction.

As we have mentioned before, the historical number of orders in one grid is supposed to have an impact on the decrease of the order matching probability of the grid. For a grid with a higher average number of orders, like the grid located around Jinrongjie, the decrease in the order matching probability is supposed to be slower, thus it is able to accept more cruising agents while still maintains a relatively decent order matching probability. For a grid with a lower number of average number of orders, like the grid located in Dongxiaoying, the decrease in the order matching probability is supposed to be faster because with a small historical number of orders, it is not able to withstand too much competition among agents. In other words, for a grid with a small number of orders, the order matching probability is decreasing rapidly when some agents are being guided into the grid.

### 4.5. Heterogeneity in reward functions

Although the determined coefficient $\alpha$ is applicable to the overall driver population at the aggregate level, it may vary from individual to individual. In other words, each driver may have a specific reward function which can be quite different from other drivers', resulting in some discrepancies in driving patterns and strategies. To examine various driving patterns of taxis drivers, especially the change in the driving patterns from the time without e-hailing to the time when e-hailing is widely adopted, Ma et al. (2019) carried out in-depth analysis using trajectories of taxi drivers in Shanghai, China, and unveiled that on average, taxi driving patterns which were previously concentrated in some central areas are now more spread out. Inspired by this, we examine the spread of a driver's trajectories using the

radius of gyration $R_g$. Radius of gyration is defined as the standard deviation of the spatial distances between a driver's location and the centroid of the driver's visited locations, i.e.,

$$R_g = \sqrt{\frac{1}{n}\sum_{i=1}^{n} r_i^2} \tag{4.1}$$

where $r_i$ is the distance between the driver's $i^{th}$ location and the centroid of the driver's visited places, and $n$ is the total number of the driver's visited locations.
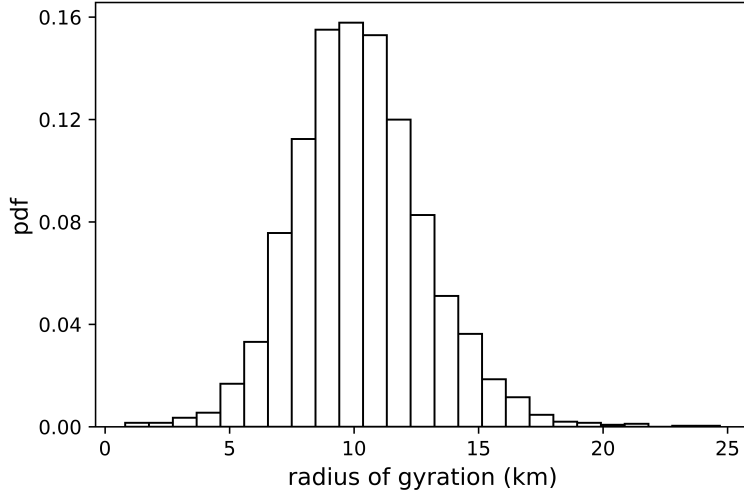


Figure 16: Distribution of radius of gyration

Figure (16) plots the distribution of the radius of gyration across the e-hailing driver population. The radius of gyration of the majority (95%) of e-hailing drivers is distributed within the range of [5.5 km, 16 km]. The e-hailing drivers on the left tail (2.5%) of the distribution has a radius of gyration below 5.5 km, and the drivers on the right tail has a radius of gyration above 16 km. The substantial discrepancy in the radius of gyration across the driver population indicates that drivers exhibit different driving patterns, which may stem from different intrinsic reward functions. For example, the trajectory of a driver on the left tail with a small radius of gyration is more concentrated in a small region, indicating that the driver may have some incentives (such as being close to home or being more familiar with the region) to stay within the region.

To shed some light on the distinction among driving patterns of e-hailing drivers, we simply take one driver on the left tail with a small radius of gyration (1.2 $km$) and one driver on the right tail with a large radius of gyration (17.8 $km$) and employ the IRL technique to disclose some underlying information regarding the driver's intrinsic reward function. Here we emphasize that understanding drivers' intrinsic reward function can be very helpful from the perspective of the platform to appropriately assign e-hailing orders. For instance, for a driver who prefers to stay within a small region, the platform is supposed to assign relatively shorter e-hailing orders to this driver. Otherwise, the driver will have a lower utility and may even be unwilling to take the assignment since the driver cares more about staying within the region over a simple higher monetary return. Thus, appropriately assigning orders can help increase the drivers' response rate and therefore the utilization of vehicle resources, as well as the passenger satisfaction, which is beneficial for all players in the e-hailing market, including the platform, the drivers, and the passengers.

Without any knowledge of the driver's demographic information, we devise the third simple reward function $\phi_3(s, s')$ as the difference between the spatial distance between state $s$ and the centroid $c$ of the driver's visited locations and the spatial distance between state $s'$ and the centroid $c$, i.e.,

$$\phi_3(s, s') = distance(s, c) - distance(s', c) \tag{4.2}$$

where $distance(s, c)$ denotes the spatial straight line distance between $s$ and $c$. The rationale of this simple reward function can be explained as follows. For a driver with a small radius of gyration, the driver

prefers to come back around the centroid after completing an order (otherwise the radius of gyration would be larger), indicating that coming back to the centroid may increase the driver's utility/intrinsic reward. $\phi_3$ exactly does this. The driver gets a positive reward when $distance(s',c) < distance(s,c)$ and a negative reward when $distance(s',c) > distance(s,c)$, meaning that it is beneficial for the driver to go back to the centroid. While for a driver with a large radius of gyration, the driver may not care about his/her distance to the centroid, and thus $\phi_3$ may not be important in the driver's intrinsic reward function. We apply the IRL technique to the observed policy of these two drivers with three simple reward functions, namely, $\phi_1$ (fare), $\phi_2$ (traveling distance), and $\phi_3$, and the derived coefficients are listed in Table (4). Again, we fix $\alpha_1 = 1$ under the assumption that the driver gets all the fare.

Table 4: Coefficients

|  | $\phi_1(s,s')$ | $\phi_2(s,s')$ | $\phi_3(s,s')$ |
|---|---|---|---|
| The driver with a small radius of gyration | 1.00 | 0.21 | 1.42 |
| The driver with a large radius of gyration | 1.00 | 0.37 | 0.17 |

The coefficient for $\phi_3$ is quite large (1.42) for the driver with a small radius of gyration and relatively small (0.17) for the driver with a large radius of gyration. This validates the effectiveness of the devised reward function $\phi_3$ in explaining the driver's intrinsic reward function when the driver's radius of gyration is small. When the driver's radius of gyration is large, the reward function $\phi_3$ does not contribute enough to the driver's underlying reward, meaning that the driver with a large radius of gyration does not gain more utility/reward by coming back to the centroid. We believe that there exist other types of simple reward functions which may be important in explaining a driver's underlying reward function when the driver has a large radius of gyration and will be left in future work.

## 5. Conclusion

Based on a large-scale real-world historical GPS traces, this paper investigated how to improve the income and rate of return of e-hailing drivers through a modified MDP approach. We proposed an MDP model which incorporates the following distinct features of drivers with e-hailing: (1) An e-hailing driver may receive a matched order before she drops off the previous passenger, thus there is no passenger seeking; (2) Different from traditional taxi that a driver has to see a passenger to find a match, e-hailing platforms very likely find a match even when the driver and the passenger are spatially far from each other. In other words, a driver's search process may end before a passenger is picked up.

We used 44,160 Didi drivers 3-day trajectories to train the model with a reward function uncovered by IRL. We then examined the optimal policy learned from our model and found that with e-hailing, the optimal policy suggests drivers to stay when they are in the city area, to move to some local areas with a high probability of receiving a request if they are in suburban areas, and to wait when they are at some places with a very high likelihood of receiving a request. To validate the effectiveness of the derived policy, a Monte Carlo simulation is conducted, and two metrics, namely the rate of return and utilization rate, are employed to compare the performance of the agent following the derived optimal policy with that of the agent following one baseline heuristic, namely, the local hotspot strategy. The comparison validates the effectiveness of the proposed model and shows that our model is able to achieve a 17.5% improvement and a 7.5% improvement over the local hotspot strategy in terms of the rate of return and the utilization rate, respectively. Also, the results show that under the guidance of the optimal policy, the agent is able to complete more order, decrease idling time, and find better orders. In addition, we disclose the reason why the agent not necessarily prefers a shorter order.

In the modified MDP model, the order matching probability captures the supply-demand ratio by its definition, considering the fact that the number of drivers in this study is sufficiently large and thus the number of unmatched orders is assumed to be negligible. Results show that the optimal policy does not necessarily guide an e-hailing driver to a grid with a high demand. Instead, the optimal policy suggests a driver to a grid with a low supply-demand ratio, i.e., a high order matching probability. To accurately capture the competition among drivers, we have devised and calibrated a dynamic adjustment strategy of the order matching probability when there are multiple drivers need recommendations. Also, the heterogeneity in reward functions across the driver population has been investigated. The devised

simple reward function $\phi_3$ is validated to explain the driving patterns of drivers with a relatively small radius of gyration and paves the way for future research on the underlying reward function of e-hailing drivers.

There are some extensions can be done to overcome some limitations of the model:

1. Although a dynamic adjustment strategy for multiple agents has been devised and validated, the model does not fully incorporate the dynamic real-time multi-agent competition. In other words, since all parameters (i.e., probabilities) are predetermined, the model is not able to fully capture the competition among agents, resulting in potential overestimation issues. A multi-agent reinforcement learning approach can thus be adopted to consider the real-time competition and can yield a more realistic optimal policy for recommendation.

2. When the grid size is small, the number of states in the MDP can be large, resulting in a higher requirement in the computation power. A hierarchical MDP model which reduces a big problem into several subproblems can be an efficient tool. For example, we can first divide the whole space into big zones and then divide each zone into finer grids.

3. A more driver-specific/personalized reward function can be investigated and is expected to incorporate more factors such as safety, stress, etc. Coupled with some prior knowledge, a more powerful IRL approach, such as the maximum entropy IRL (Ziebart et al., 2008), could be an promising way to uncover more information of the underlying reward function that results in the demonstrated behavior.

## Acknowledgments

## References

Alemi, F., Circella, G., Mokhtarian, P., Handy, S., May 2019. What drives the use of ridehailing in California? Ordered probit models of the usage frequency of Uber and Lyft. Transportation Research Part C: Emerging Technologies 102, 233–248.

Battifarano, M., Qian, Z. S., Oct. 2019. Predicting real-time surge pricing of ride-sourcing companies. Transportation Research Part C: Emerging Technologies 107, 444–462.

Bertsekas, D. P., 2000. Dynamic Programming and Optimal Control, 2nd Edition. Athena Scientific.

Di, X., Ban, X. J., Nov. 2019. A unified equilibrium framework of new shared mobility systems. Transportation Research Part B: Methodological 129, 50–78.

Di, X., Liu, H. X., Davis, G. A., 2010. Hybrid extended kalman filtering approach for traffic density estimation along signalized arterials: Use of global positioning system data. Transportation Research Record 2188 (1), 165–173.

Di, X., Liu, H. X., Zhu, S., Levinson, D. M., 2017. Indifference bands for boundedly rational route switching. Transportation 44, 11691194.

Gao, Y., Jiang, D., Xu, Y., Aug. 2018. Optimize taxi driving strategies based on reinforcement learning. International Journal of Geographical Information Science 32 (8), 1677–1696.

Ge, Y., Xiong, H., Tuzhilin, A., Xiao, K., Gruteser, M., Pazzani, M., 2010. An Energy-efficient Mobile Recommender System. In: Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '10. ACM, New York, NY, USA, pp. 899–908.

He, F., Shen, Z.-J. M., Sep. 2015. Modeling taxi services with smartphone-based e-hailing applications. Transportation Research Part C: Emerging Technologies 58, 93–106.

He, F., Wang, X., Lin, X., Tang, X., 2018. Pricing and penalty/compensation strategies of a taxi-hailing platform. Transportation Research Part C: Emerging Technologies 86, 263–279.

Hu, X., Gao, S., Chiu, Y.-C., Lin, D.-Y., Jan. 2012. Modeling routing behavior for vacant taxicabs in urban traffic networks. Transportation Research Record 2284 (1), 81–88.

Hwang, R.-H., Hsueh, Y.-L., Chen, Y.-T., Sep. 2015. An effective taxi recommender system based on a spatio-temporal factor analysis model. Information Sciences 314, 28–40.

Li, B., Zhang, D., Sun, L., Chen, C., Li, S., Qi, G., Yang, Q., Mar. 2011. Hunting or waiting? Discovering passenger-finding strategies from a large-scale real-world taxi dataset. In: 2011 IEEE International Conference on Pervasive Computing and Communications Workshops (PERCOM Workshops). pp. 63–68.

Lin, K., Zhao, R., Xu, Z., Zhou, J., 2018. Efficient Large-Scale Fleet Management via Multi-Agent Deep Reinforcement Learning. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. KDD '18. ACM, New York, NY, USA, pp. 1774–1783.

Liu, L., Andris, C., Ratti, C., Nov. 2010. Uncovering cabdrivers' behavior patterns from their digital traces. Computers, Environment and Urban Systems 34 (6), 541–548.

Liu, S., Araujo, M., Brunskill, E., Rossetti, R., Barros, J., Krishnan, R., Jun. 2013. Understanding Sequential Decisions via Inverse Reinforcement Learning. In: 2013 IEEE 14th International Conference on Mobile Data Management. Vol. 1. pp. 177–186.

Ma, Q., Yang, H., Zhang, H., Xie, K., Wang, Z., Oct. 2019. Modeling and Analysis of Daily Driving Patterns of Taxis in Reshuffled Ride-Hailing Service Market. Journal of Transportation Engineering, Part A: Systems 145 (10), 04019045.

Markou, I., Kaiser, K., Pereira, F. C., May 2019. Predicting taxi demand hotspots using automated Internet Search Queries. Transportation Research Part C: Emerging Technologies 102, 73–86.

Moreira-Matias, L., Gama, J., Ferreira, M., Damas, L., Sep. 2012. A predictive model for the passenger demand on a taxi network. In: 2012 15th International IEEE Conference on Intelligent Transportation Systems. pp. 1014–1019.

Ng, A. Y., Russell, S. J., 2000. Algorithms for Inverse Reinforcement Learning. In: Proceedings of the Seventeenth International Conference on Machine Learning. ICML '00. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 663–670.

Phithakkitnukoon, S., Veloso, M., Bento, C., Biderman, A., Ratti, C., Nov. 2010. Taxi-Aware Map: Identifying and Predicting Vacant Taxis in the City. In: Ambient Intelligence. Lecture Notes in Computer Science. Springer, Berlin, Heidelberg, pp. 86–95.

Powell, J. W., Huang, Y., Bastani, F., Ji, M., 2011. Towards Reducing Taxicab Cruising Time Using Spatio-temporal Profitability Maps. In: Proceedings of the 12th International Conference on Advances in Spatial and Temporal Databases. SSTD'11. Springer-Verlag, Berlin, Heidelberg, pp. 242–260.

Puterman, M. L., 1994. Markov Decision Processes: Discrete Stochastic Dynamic Programming, 1st Edition. John Wiley & Sons, Inc., New York, NY, USA.

Qian, X., Ukkusuri, S. V., Jun. 2017. Taxi market equilibrium with third-party hailing service. Transportation Research Part B: Methodological 100, 43–63.

Qu, M., Zhu, H., Liu, J., Liu, G., Xiong, H., 2014. A Cost-effective Recommender System for Taxi Drivers. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '14. ACM, New York, NY, USA, pp. 45–54.

Rong, H., Zhou, X., Yang, C., Shafiq, Z., Liu, A., 2016. The Rich and the Poor: A Markov Decision Process Approach to Optimizing Taxi Driver Revenue Efficiency. In: Proceedings of the 25th ACM International on Conference on Information and Knowledge Management. CIKM '16. ACM, New York, NY, USA, pp. 2329–2334.

Shou, Z., Di, X., Nov. 2018. Similarity analysis of frequent sequential activity pattern mining. Transportation Research Part C: Emerging Technologies 96, 122–143.

Sirisoma, R., Wong, S. C., Lam, W. H. K., Wang, D., Yang, H., Zhang, P., 2010. Empirical evidence for taxi customer-search model. Proceedings Of The Institution Of Civil Engineers: Transport 163 (4), 203–210.

Sutton, R. S., Barto, A. G., 1998. Introduction to Reinforcement Learning, 1st Edition. MIT Press, Cambridge, MA, USA.

Szeto, W. Y., Wong, R. C. P., Wong, S. C., Yang, H., Jul. 2013. A time-dependent logit-based taxi customer-search model. International Journal of Urban Sciences 17 (2), 184–198.

Tan, X., Zhang, J., Cao, W., Li, J., Zheng, Y., 2018. When will you arrive? estimating travel time based on deep neural networks. In: AAAI.

Tang, J., Jiang, H., Li, Z., Li, M., Liu, F., Wang, Y., 2016. A two-layer model for taxi customer searching behaviors using gps trajectory data. IEEE Transactions on Intelligent Transportation Systems 17 (11), 3318–3324.

Verma, T., Varakantham, P., Kraus, S., Lau, H. C., Jun. 2017. Augmenting decisions of taxi drivers through reinforcement learning for improving revenues. Proceedings of the Twenty-Seventh International Conference on Automated Planning and Scheduling ICAPS 2017: Pittsburgh, June 18-23, 409–417.

Wang, Y., Gu, J., Wang, S., Wang, J., Aug. 2019. Understanding consumers willingness to use ride-sharing services: The roles of perceived value and perceived risk. Transportation Research Part C: Emerging Technologies 105, 504–519.

Wong, K. I., Wong, S. C., Bell, M. G. H., Yang, H., 2005. Modeling the bilateral micro-searching behavior for urban taxi services using the absorbing markov chain approach. Journal of Advanced Transportation 39 (1), 81–104.

Wong, K. I., Wong, S. C., Yang, H., Nov. 2001. Modeling urban taxi services in congested road networks with elastic demand. Transportation Research Part B: Methodological 35 (9), 819–842.

Wong, K. I., Wong, S. C., Yang, H., Wu, J. H., Dec. 2008. Modeling urban taxi services with multiple user classes and vehicle modes. Transportation Research Part B: Methodological 42 (10), 985–1007.

Wong, R. C. P., Szeto, W. Y., Wong, S. C., Nov. 2014a. A cell-based logit-opportunity taxi customer-search model. Transportation Research Part C: Emerging Technologies 48, 84–96.

Wong, R. C. P., Szeto, W. Y., Wong, S. C., 2015a. Sequential Logit Approach to Modeling the Customer-Search Decisions of Taxi Drivers. Asian Transport Studies 3 (4), 398–415.

Wong, R. C. P., Szeto, W. Y., Wong, S. C., Oct. 2015b. A two-stage approach to modeling vacant taxi movements. Transportation Research Part C: Emerging Technologies 59, 147–163.

Wong, R. C. P., Szeto, W. Y., Wong, S. C., Yang, H., Jan. 2014b. Modelling multi-period customer-searching behaviour of taxi drivers. Transportmetrica B: Transport Dynamics 2 (1), 40–59.

Wong, S., Yang, H., Jan. 1998. Network Model of Urban Taxi Services: Improved Algorithm. Transportation Research Record: Journal of the Transportation Research Board 1623, 27–30.

Xu, Z., Li, Z., Guan, Q., Zhang, D., Li, Q., Nan, J., Liu, C., Bian, W., Ye, J., 2018. 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. KDD '18. ACM, New York, NY, USA, pp. 905–913.

Yang, H., Fung, C. S., Wong, K. I., Wong, S. C., Jun. 2010a. Nonlinear pricing of taxi services. Transportation Research Part A: Policy and Practice 44 (5), 337–348.

Yang, H., Leung, C. W. Y., Wong, S. C., Bell, M. G. H., Sep. 2010b. Equilibria of bilateral taxicustomer searching and meeting on networks. Transportation Research Part B: Methodological 44 (8), 1067–1083.

Yang, H., Wong, S. C., May 1998. A network model of urban taxi services. Transportation Research Part B: Methodological 32 (4), 235–246.

Yang, H., Wong, S. C., Wong, K. I., Nov. 2002. Demand-supply equilibrium of taxi services in a network under competition and regulation. Transportation Research Part B: Methodological 36 (9), 799–819.

Yang, H., Yang, T., May 2011. Equilibrium properties of taxi markets with search frictions. Transportation Research Part B: Methodological 45 (4), 696–713.

Yang, T., Yang, H., Wong, S. C., Oct. 2014. Taxi services with search frictions and congestion externalities. Journal of Advanced Transportation 48 (6), 575–587.

Yu, X., Gao, S., Hu, X., Park, H., Mar. 2019. A Markov decision process approach to vacant taxi routing with e-hailing. Transportation Research Part B: Methodological 121, 114–134.

Yuan, J., Zheng, Y., Zhang, L., Xie, X., Sun, G., 2011. Where to Find My Next Passenger. In: Proceedings of the 13th International Conference on Ubiquitous Computing. UbiComp '11. ACM, New York, NY, USA, pp. 109–118.

Zhang, K., Chen, Y., Nie, Y. M., Dec. 2019a. Hunting image: Taxi search strategy recognition using Sparse Subspace Clustering. Transportation Research Part C: Emerging Technologies 109, 250–266.

Zhang, K., Jia, N., Zheng, L., Liu, Z., Nov. 2019b. A novel generative adversarial network for estimation of trip travel time distribution with trajectory data. Transportation Research Part C: Emerging Technologies 108, 223–244.

Zhou, X., Rong, H., Yang, C., Zhang, Q., Khezerlou, A. V., Zheng, H., Shafiq, M. Z., Liu, A. X., 2018. Optimizing Taxi Driver Profit Efficiency: A Spatial Network-based Markov Decision Process Approach. IEEE Transactions on Big Data, 1–1.

Ziebart, B. D., Maas, A., Bagnell, J. A., Dey, A. K., 2008. Maximum Entropy Inverse Reinforcement Learning. In: Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3. AAAI'08. AAAI Press, Chicago, Illinois, pp. 1433–1438.