# Spatio-temporal Incentives Optimization for Ride-hailing Services with Offline Deep Reinforcement Learning

YANQIU WU, New York University, USA

QINGYANG LI, DiDi, China

ZHIWEI QIN, DiDi, China

A fundamental question in any peer-to-peer ride-sharing system is how to, both effectively and efficiently, meet the request of passengers to balance the supply and demand in real time. On the passenger side, traditional approaches focus on pricing strategies by increasing the probability of users' call to adjust the distribution of demand. However, previous methods do not take into account the impact of changes in strategy on future supply and demand changes, which means drivers are repositioned to different destinations due to passengers' calls, which will affect the driver's income for a period of time in the future. Motivated by this observation, we make an attempt to optimize the distribution of demand to handle this problem by learning the long-term spatio-temporal values as a guideline for pricing strategy. In this study, we propose an offline deep reinforcement learning based method focusing on the demand side to improve the utilization of transportation resources and customer satisfaction. We adopt a spatio-temporal learning method to learn the value of different time and location, then incentivize the ride requests of passengers to adjust the distribution of demand to balance the supply and demand in the system. In particular, we model the problem as a Markov Decision Process (MDP). The problem is solved in two steps: 1) based on historical trip data from the ride-hailing platform, we propose a deep reinforcement learning based method with constrained actions to summarize demand and supply patterns into a Spatio-Temporal network, 2) to solve the budget constraint problem in the spatio-temporal incentive, we formulate an integer programming step in real-time policy learning process, where each state-action pair is valued in terms of immediate reward, future gains and discount budget. Through experiments, we demonstrate that our proposed approach can deliver improvement on the marketplace efficiency.

CCS Concepts: • **Computing methodologies → Planning for deterministic actions**; **Q-learning**.

Additional Key Words and Phrases: reinforcement learning, spatio-temporal incentivization, off-line learning, markov decision processes

## 1 INTRODUCTION

With rapid development of the Global Position System(GPS), and on-demand ride-hailing services such as Uber, Lyft and Didi Chuxing, significant improvements has been achieved over traditional taxi systems. Massive amount of trip data become available, offering much more opportunities for providing more intelligent and convenient services and

Authors' addresses: Yanqiu Wu, New York University, USA, yanqiu.wu@nyu.edu; Qingyang Li, DiDi, China, qingyangli@didiglobal.com; Zhiwei Qin, DiDi, China, qinzhiwei@didiglobal.com.

a surge in passion in research fields such as driving routing planning [12][22], order dispatching [6][20][11], supply chain management[18], driver program [16][15] and demand prediction [25][1].

One of the most important task in online ride-hailing platforms is to balance supply and demand situation in advance. Once passenger's request has been answered, the trip would bring both the passenger and the driver to the destination, which can be viewed as a process of relocating the driver to the destination. The driver can answer future requests at the new location. For example, when a ride order travels to a hot (high demand) zone where drivers are more needed, after fulfilling the ride, the driver is more likely to receive new requests there. If the situation of traversing to a hotter zone persists, the driver can receive more gain in the long turn. Hence, we see that the value of a trip goes beyond the trip fare and depends also on the spatio-temporal long-term values.

Dynamic pricing and incentives are one of the levers in the ride-hailing system that can influence both supply and demand distributions to make them more aligned, and hence, improving customer experience through higher driver income and request answer rate. Dynamic pricing and spatial pricing are mainly explored by previous researchers. Some work focused on approximate dynamic programming(ADP) [24] and model predictive control(MPC) [14]. Other researchers build up spatial equilibrium [2, 7, 26]. These research studies have some limitations. For example, they can examine dynamic and spatial pricing separately but not simultaneously. Recently, to achieve spatio-temporal incentivization, some work of reinforcement learning enhanced agent-based modeling and simulation system have been explored [3]. However, to optimize the reinforcement learning policy, the method requires to build up a data-driven simulation environment to mimic a real-world ride-sourcing platform's operations on a city network to enable online training.

Recently, offline reinforcement learning [4][10][9] has drawn a lot of attention since it promises to learn effective policies from previously-collected, static datasets without further interaction. Because offline RL does not rely on the built environment, it has achieved great success in many real-world application scenarios, such as recommendation system, real-time decision-making system, advertising optimization, robotics and so on. BCQ [5] proposed a batch-constrained reinforcement learning method close to on-policy with respect to the available data to solve the extrapolation error problem. BEAR [8] is a practical solution to reduce the bootstrapping error accumulation.

However, our problem can not be solved directly with the existing offline RL method. In the real ride-hailing scenarios, when we model the spatio-temporal incentives optimization problem as a Markov Decision Process(MDP), incentive is treated as the action and it will bring the related cost when optimizing the policy. There are some related work [21] to study the MDP problem under budget constraints while none of them will work in the offline RL scenarios. It brings more challenges when we constrain the budget of the whole city from a global view, the solution space of this problem will be huge and it becomes a NP-hard problem.

Our work adopts reinforcement learning methods, but instead of training a simulated ride-sourcing platform and an on-policy agent, we formulate our problem in the offline reinforcement learning setting, and utilize previously collected data, without additional online interactions with the environment. The goal for spatio-temporal incentive is not only to fulfill the ride requests of the current passengers, but also to optimize the anticipated future gains. In this study, we propose an offline deep reinforcement learning to learn the spatio-temporal value network to model the demand and mobility patterns for the whole city. In addition, to solve the budget constraint problem during the policy learning, we formulate it as an integer programming process to conduct the policy optimization within the limitation of budget threshold and learn the state-action network simultaneously. To verify the effectiveness of the proposed method, we apply it to a real ride-hailing system. Through comparative evaluations, our proposed method shows significant improvements in the real-world application.
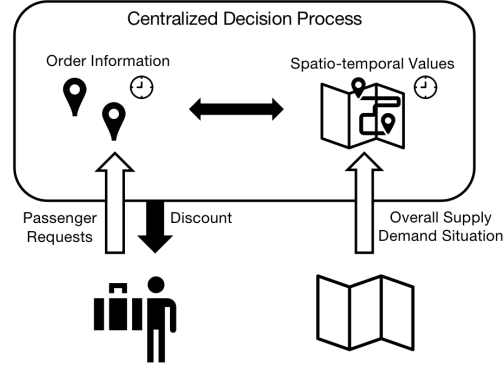
Fig. 1. Overview of spatio-temporal incentives optimization in ride-hailing services.

The contribution of this work can be summarized as follows:

- We propose a novel offline reinforcement learning to dynamically adjust the spatial-temporal distribution of demand to balance driver supply and customer demands. To the best of our knowledge, this is the first work of offline RL to optimize spatio-temporal problem on the demand side.
- We tackle the budget constraint problem in the offline RL setting by formulating the constraint RL problem as an integer programming process to learn the spatio-temporal value network as well as control the resource consumption of policy.
- Experiments are conducted under different times of the week(weekdays, weekends) and from multiple aspects to analyze supply and demand, including comparing to SOTA offline RL algorithms. We largely improve both revenue and efficiency in balancing supply and demand across the city.

In what follows, we will describe the passenger pricing problem in section 2 and formulate our MDP in section 3. The detailed explanation of our deep Q-network approach combined with Integer Programming can be found in section 4. The learning and optimization capabilities of our deep reinforcement learning algorithm and the advantage of the integer programming method through a set of experiments in section 5 using real data from a major ride-hailing platform. We close the paper with a few concluding remarks in section 6.

## 2  RELATED WORK

### 2.1  Background and System Overview of Spatio-temporal Incentive

Recently, on-demand online ride-hailing platform like Uber, Lyft, and DiDi provide a more efficient way of transportation. On the passengers' side, before they actually request a ride order, they need to provide information of the ride, including the starting location, destination and the time of the ride order. We name this stage as ride inquiry. Then, these online ride services have a centralized system to give some information back to the passengers, such as estimated travel time and estimated ride fee. Finally, based on these information, passengers can decide whether they actually want to place the ride order. We name this stage as ride request. The conversion probability from ride inquiry to ride request is defined as estimated-call-rate (ECR). Naturally, while other information of the ride keeps unchanged, the cheaper the

fee is, the higher the ecr would be. Once the passenger requests for a ride, the central dispatching system will try to match a driver to pick up the passenger to finish the ride. The probability of a ride being successfully completed by a driver is called completion-rate(CR). CR will not be affected by ride pricing.

The optimization of passenger pricing and incentive is crucial to balance the demand and supply since the incentive strategy will directly have influences on the distribution of demand. One special case is that when the amount of supply is much larger than the demand, which means there are more idle drivers in this region, the platform should give more discount to passengers to increase the demand amount, helping drivers complete more orders and increase drivers' income. Another case is that when amount of supply is less than the demand, which means there are a lot of trip request in this region. However, different order requests have different destinations, which will take the driver to a different geographic location in the next moment and affect the driver's future income. If the driver goes from the hot zone to the cold zone, it will reduce the driver's income from the next order, otherwise it will increase the driver's probability of taking orders and increase income. Therefore, the optimization of the pricing strategy is crucial to the driver's income, the balance of supply and demand, and the improvement of the overall trading market efficiency.

The incentives system is responsible for distributing real-time discounts to passengers with open trip inquiries. Figure 1 shows a system overview of how it works in practice.Every time a passenger enters the ride inquiry stage, a centralized decision system would decide on the discount rate of the current ride. Initially, the decision only affects passenger's ecr of the order. Once the passenger enters the ride request stage and is successfully picked up by a driver, the discount would only takes effect at the payment stage.

Given both historical ride data and real-time order information, we can potentially optimize the pricing system's global efficiency in a long horizon by spatio-temporal pricing. When a passenger requests a ride and the order is finished, he or she has moved from the starting location to the destination. Meanwhile, the driver who carries out the ride, is also relocated to the destination. Hence, a trip changes the spatio-temporal state of a driver, which together with demand distributions, forms the basis for future order matching and vehicle repositions. The spatio-temporal state values are the quantities that characterize these long-term effects and can thus guide the incentives decisions to optimize global metrics, e.g., total driver income.

## 3  LEARNING

The learning step aims to provide a quantitative measure of the spatio-temporal patterns of the ride supply and passenger demand across the whole city. Given historical off-line ride data, we build a Markov Decision Process (MDP) representing rides happening on the platform. The learned value functions provide quantitative values for each spatio-temporal state, which is later used in the discount planning step.

### 3.1  Problem Statement

*MDP definition.* The MDP we build here is in a global view. In [23], it proposed a tabular based reinforcement learning method to solve the problem of order dispatch in the ride-hailing platform from the perspective of local view, which aims to match drivers with different ride orders, each driver is modeled as an agent. However, in this study, unlike drivers, if a passenger requests a ride from location A to location B, there's no guarantee that he or she will request the next ride starting from location B. Hence, we can not model each independent passenger as an agent. Instead, we adopt the global-view setting which we do not distinguish individual passengers. Each passenger is regarded as the same. For example, if one ride transits from A at time t and arrives at B at time t+1, any ride leaving B at time t+1 can be regarded as our possible next transit.
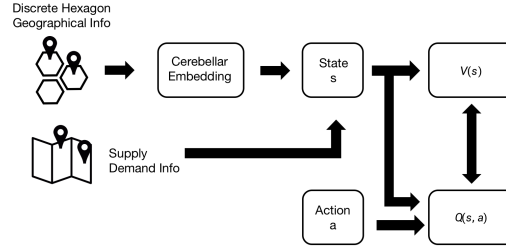
Fig. 2. Architecture of the proposed algorithm

When a passenger requests a ride, we have the information of the starting location, destination, time of the request, estimated travel time and estimated ride cost. Our MDP formulation uses these available information. A trip transition consists of order request, pick-up and completion: The passenger enters the origin and destination location. Ride-hailing platforms provide estimated travel time and fee. The trip moves the driver and passenger to the destination. The driver earns an immediate reward (trip fee) from this transition. We list the key elements of our MDP formulation below.

*State.* The state is defined to indicate geo-location information of the starting location and the time of the ride request. For simplicity, we quantize time periods into semi-hours. Formally, we define $s = (g, t, f) \in S$, where $g \in G$ represents the grid region where the trip starts, $t \in T$ is the time index. We also differentiate the time for weekday and weekend. And $f \in F$ contains the additional contextual features, such as the statistics of supply and demand of the current grid $g$.

*Action.* The action is defined as the discount assigned to the trip. Action space is discrete: 0.75, 0.8, 0.85, 0.9, 0.95, 1, corresponding to 25% off, 20% off, 15% off, 10% off, 5% off and no discount of the ride fee. Naturally, as the discount goes up of the ride fee, the probability of a ride inquiry becoming a actual ride request would increase.

*Reward.* In order dispatch setting, reward is defined as the total fee collected for the trip. Intuitively, in our problem, the reward can be defined as: $r(s_t, a_t) = \delta(ecr_t) * GMV_t$, where $\delta(ecr_t)$ represents the change of ecr[1] under discount $a_t$ and $GMV_t$ represents the trip fare of trip $s_t$. The optimal value network of a state s, $V^*(s)$, would be the optimal estimated future delta Gross Merchandise Volume(GMV).

*Episode.* An episode is one complete day, from 0:00 am to 23:59pm. Trips crosses midnight are regarded as terminal state.

*State transition* and *reward distribution.* The discount distribution have to be made in a rolling horizon manner. In real work application, at the time of a ride request, the actual trip fare and travel time are unknown. Our system have to depends on the estimated arrival time and estimated price to make decisions. The error distribution of these estimations are implicitly captured by state transitions and reward distributions.

*Discount factor.* The discount factor in the MDP definition controls how far we look in to the future. In the order dispatch problem, the rewards defined in [23] are also discounted. Here in our passenger discount application, we adopt the same discounting methods. For example, for an order which takes $T$ time slots(semi hours) with trip fare $GMV$ and a discount factor $\gamma$, the final reward is given by:

$$R_\gamma = \delta(ecr) * \left( \sum_{t=0}^{T-1} \gamma^t \frac{R}{T} \right) \tag{1}$$

---

[1]ECR represents the probability of a user calling an order, which ranges from 0 to 1.

For example, suppose a passenger requests a order from A to B at time 00:00. The trip is estimated to be completed in 40 minutes and costs 40 dollars. Suppose we assigns 10% off discount to the order, which increase the passenger's ecr by 15%. We use a discount factor of $\gamma = 0.9$ and segment the time slots into 10-minutes windows. In our model, this order will transit from $(A, 0, f_a)$ to $(B, 4, f_b)$ with action $a = 0.9$ and $r = \delta(ecr) * (10 + 10 * 0.9 + 10 * 0.9^2 + 10 * 0.9^3) = 0.15 * 34.39 = 5.1585$.

## 4 ALGORITHM

In this section, we proposed an offline deep reinforcement learning method to learn the long-term spatio-temporal value for the whole city, to dynamically change the distribution of demand and achieve a state of balance between supply and demand, improving the efficiency of the marketplace. For more detailed information of our proposed method, please refer to the pseudocode provided in Algorithm 1 and the architecture overview in figure 2.

### 4.1 Deep Q-network

To solve the MDP formulation, we adopt the model-free approach and use the deep reinforcement learning framework proposed in [13]. In the deep Q-network framework, the mini-batch update is a step for solving a bootstrapped regression model with the Mean Squared Error loss function:

$$\left( Q(s_t, a_t | \theta) - \left( r(s_t, a_t) + \gamma \arg\max_{a_{t+1} \in A} Q(s_{t+1}, a_{t+1} | \theta') \right) \right) \tag{2}$$

where $\theta'$ is the weights for the target Q network, $A$ is the action space. To improve training stability and avoid overestimation, we use Double Q learning [19]. Specifically, in addition to the Q network, we use a corresponding Q target network. The targets in equation (2) is modified so that the argmax is evaluated by Q network, and the highest valued action is used by the target Q network as follows:

$$r(s_t, a_t) + \gamma Q_{targ}(s_t, \arg\max_a Q(s_t, a)) \tag{3}$$

### 4.2 Model Training

Building a realistic simulator to imitate city-scale driver and passenger dynamics is very challenging. However, a large number of historical trips are usually available. Hence, we use historical trip data for training, which means our training is off-line and we do not interact with any environment during training. The historical data is generated by some existing behavior policy. Each trip $x$ defines a transition $(s_t, a_t, r_t, s_{t+1})$. These data are retrieved from data warehouse and store in a replay buffer for training similar to [20]. During training, each mini-batch is sampled from the replay buffer to train the Q and V spatio-temporal network.

### 4.3 Cerebellar embedding

Having a good state representation is usually the key step to solving a practical problem with neural networks. In our spatio-temporal incentives optimization problem, we require the parse of complicated state information as the basis for long term reasoning. Hence, we adopt the cerebellar embedding scheme used in order dispatching problem[17]. Cerebellar embedding combines CMACs with embedding to obtain a distributed state representation[17] that is generalizable, extensible and robust. CMACs involves multiple overlapping tilings of the state space. The total number of tiles is the size of the conceptual memory. The mapping from input to tiles is done such that input points close

---

**Algorithm 1** Spatio-Temporal Network

---

Load historical data into the replay buffer $\mathcal{M}$, state $s = (l, t, f)$
Ranomly initialize the spatio-temproal network $\mathcal{V}$ with weights $\phi$ and state-action value network $Q$ with weights $\theta$
Initialize the target network $\hat{\mathcal{V}}$ and $\hat{Q}$ with the same weights as $\mathcal{V}$ and $Q$
**for** t = 1,2, ..., T **do**
   Sample a random mini-batch $\mathcal{B} = \{(s_j, a_j, r_j, s_{j+1})\}$ from $\mathcal{M}$.
   Calculate $\mathcal{V}$ target $v_j$ with
   $v_j = \hat{Q}_{\theta'}(s_j, \arg\max_a Q_\theta(s_j, a))$
   Calculate $Q$ target $y_j$ with
   $y_j = r_j + \gamma(1 - d)\hat{\mathcal{V}}_{\phi'}(s_{j+1})$ {d is the terminal state signal}
   Update $Q$ by one step of gradient descent on
   $\nabla_\theta \frac{1}{|\mathcal{B}|} \sum (Q_\theta(s_j, a_j) - y_j)^2$
   Update $\mathcal{V}$ by one step of gradient descent on
   $\nabla_\phi \frac{1}{|\mathcal{B}|} \sum (\mathcal{V}_\phi(s_j) - v_j)^2$
   Calculate advantage matrix $\mathcal{A}$ as in equation 5
   Calculate cost matrix $C$ as in equation 6
   Solve policy distribution $X$ using $A$ and $C$ as in equation 4.
   Add $\mathcal{B}' = \{(s_j, a'_j, r'_j, s'_{j+1})\}$ in $\mathcal{M}$ where $a'_j$ are given by $X$.
   **if** t mod update frequency = 0 **then**
     Update the target network with:
     $\theta' \leftarrow \theta$
     $\phi' \leftarrow \phi$
   **end if**
**end for**
return $\hat{\mathcal{V}}$ and $\hat{Q}$

---

together in the input space have considerable overlap between their set of activated tiles. The output of CMACs is computed as the sum of the weights of the activated tiles.

The cerebellar embedding extends CMACs. An randomly initialized embedding matrix is used as the actual memory and the mapping is implemented using a sparse representation. It defines multiple tiling functions and each function would map the continuous input to a unique string id indicating one discretized region of the state space. The full details of cerebellar embedding can be found in [17].

To quantize the geographical space, we also adopt the hierarchical coarse-coding in the location space as in [17]. Hexagon is used as the tile shape because hexagons have only one distance between a hexagon center-point and its neighbors.

### 4.4 Extrapolation error

Off-line reinforcement learning has the problem of extrapolation error [8]. Exploration error happens due to mismatch in the distribution of data produced by the policy and the distribution of data contained in the batch. Fujimoto et al. [5] introduce batch-constrained reinforcement learning, where agents are trained to maximize accumulated rewards and minimize the distance of selected actions to the data in the batch. Recall that our action space is discrete. Hence, to tackle extrapolation error, we perform a simple action search combined with Q networks, so that selected highest valued action has been seen in the batch. $\hat{A}(s) = \{a | (g(s), a) \in B\}$ where $g(s)$ is the discretized spatio-temporal grid

(a) GMV comparison in terms of calling orders          (b) GMV comparison in terms of finished orders
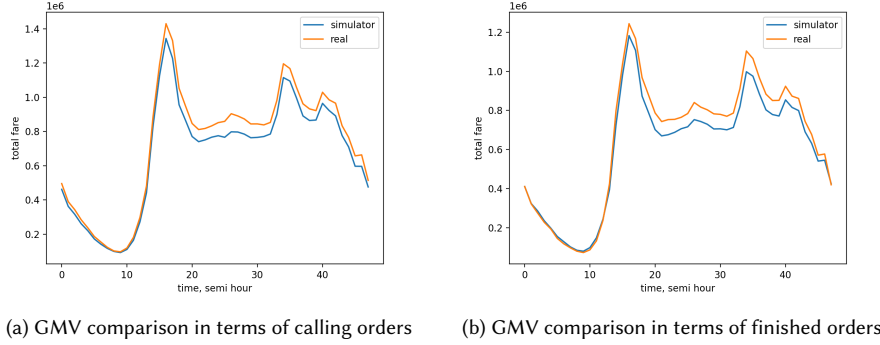
Fig. 3. The simulator calibration in terms of GMV. The red curves plot the GMV values of real data averaged in 24 hours of one day, in 30-minute time granularity. The blue curves are simulated results averaged over 48 semi-hours.

that state s falls into and $B$ is replay buffer. If the action search returns an empty action set, we would randomly select an action instead of the highest valued action to ensure exploration.

## 4.5 Terminal state

Naturally, we would define the trips across the day as our terminal states. However, in reality, especially in big cities, there are still a large amount of trips happening around midnight. Hence, in order to be consistent with real world situations, we define trips crossing 3:00am as our terminal states.

## 4.6 Integer Programming

Recall that our action space is discrete, given the deep Q-network framework, a natural policy is to select actions that maximize the learned Q values. However, in this specific application environment, ride-hailing platforms have a total budget for distributing discounts to passengers. Hence, we need to make sure that discount distribution induced by our policy will not exceed the total budget. Simply selecting the highest valued action does not guarantee this constrain. Instead, we form the optimal discount distribution problem as an Integer Programming problem:

$$\text{maximize } A \circ X$$
$$\text{subject to } C \circ X \leqslant B, X \geqslant 0 \tag{4}$$

where $A$ represents the value matrix of the trips we want to solve, $C$ represents the corresponding cost for discounts, $B$ is the total budget.

Formally, matrix $X$ is a binary matrix such that for each trip, only one of the 6 actions has the corresponding value of 1. We defined our value matrix $A$ as a N by 6 matrix to represent the different values each trip can get under different actions. N represents the number of trips. $\beta$ is a hyper-parameter used to weigh between the trip fare and the Spatio-Temporal value transitions traveling from the starting location to the destination. cr[2] indicates the probability of successful completion of the order.

$$A = \delta \text{ecr} * \text{cr} * (\beta \text{GMV} + (1 - \beta)\gamma(1 - d)\mathcal{V}(s') - \mathcal{V}(s)) \tag{5}$$

---

[2]cr denotes the probability that a calling order is accepted by the driver and the order is finally completed. Its value ranges from 0 to 1.

Manuscript submitted to ACM

The cost matrix $C$ is similarly defined as follows:

$$C = (1 - a) * \text{GMV} \tag{6}$$

where $a$ represents the discrete actions which take values of 0.75, 0.8, 0.85, 0.9, 0.95 and 1.

By solving the integer programming problem, we receive a discount distribution that optimize for the combination of trip fares and Spatio-Temporal value transitions, meanwhile, we do not exceed budget.

## 5  EXPERIMENT

In this section, we will discuss the experiment settings and results. We use historical ExpressCar trip data obtained from a major ride-hailing platform as our training data. The data set consists of over three millions of trips happening from April 13th to April 26th, 2020 and is divided into training set (one week) and testing set (one week). Data points in each mini-batch are viewed as samples on the Q value function. We used a discount factor $\gamma = 0.9$ and normalized all state vectors with their population mean and standard deviation. We found the pre-processing was necessary for a stable training.

To provide a more complete and direct understanding of the effectiveness of the proposed algorithm, we evaluate it from three different aspects, including a environment simulation, algorithm performance comparison and map visualization.

### 5.1  Environment Simulation



(a) Overall discount distribution

(b) Our policy discount distribution with total fare

(c) Historical discount distribution with total fare

Fig. 4. Action distribution comparison

In reality, whether a passenger go from ride inquiry stage to actual ride request is either yes or no, however, boolean variable is not convenient for us to model the changes in the customers' call-rate. Hence, based on historical data, we build a model to simulate the transitions and turn the boolean variable into ECR values. ECR ranges from 0 to 1 exclusive of boundaries. Before using estimated probability in our DRL algorithm training, we first evaluate the performance compared to historical data.

In Figure 3a, for rides happening in each 30-minutes interval in a consecutive 14 days, we sum up the ride fares. The x axis represents each semi-hour and ranges from 1 to 48. The y axis is the sum of all the fares during the 14 days, we name it as Gross Merchandise Volume (GMV). In real world situation, the expected revenue contribution of each trip inquiry equals its ride fee times the ecr_label, which is the boolean variable of whether a passenger turns from ride inquiry to ride request, and then times the applied discount. The results are plotted in orange and labeled as real. Meanwhile, for the same number of rides, we sum up the total simulation fare. The simulation fare of each trip is

Table 1. Total trip fare comparison

| Algorithms | Historical | IP Policy(ours) | BCQ | BEAR | BCQ+IP | BEAR+IP |
|---|---|---|---|---|---|---|
| Budget(weekdays) | 471489.1 | 471489.07 | 320090.79 | 330953.68 | 471489.07 | 471489.06 |
| Total trip fare(weekdays) | 9935198.67 | 10304260.73 | 9839823.21 | 9850487.36 | 10302276.97 | 10238700.15 |
| Increase(weekdays) | - | 3.7% | -0.96% | - 0.85% | 3.7% | 3.1% |
| Budget(weekends) | 132813.58 | 132813.58 | 109123.81 | 115062.95 | 132813.57 | 132813.57 |
| Total trip fare(weekends) | 4205931.14 | 4339827.87 | 4178529.71 | 4182857.98 | 4339146.09 | 4314800.63 |
| Increase(weekends) | - | 3.18% | -0.65% | -0.55% | 3.17% | 2.59% |

calculated as the product of the ride fee, ecr, and the applied discount of each trip. Ecr is our modeled probability of the passenger turning from ride inquiry to ride request under the applied discount of each trip. The results are plotted in blue and labeled as simulation. We can find that the simulated GMV trend is basically consistent with the real historical GMV trend. There will be slight differences in some periods, but the overall simulation error is still relatively low.

In real world application, when a passenger decide not to request a ride order after inquiring the price, the gain of the trip is zero. The same trip in our simulation, would have a gain larger than zero, because the model probability is always larger than zero. Reversely, if a passenger decide to place a request, the ecr_label is one. In this case, the product equals the ride fee times the applied discount. Under the same situation, since our modeled probability, ecr, is less than one, the fare in simulation is smaller than the fare in reality. Ideally, when we have a very large numbers of trips, the sum of their total fares would be very close under two different modes by Law of Large Numbers. Based on Figure 3a, we can discover that, during the 48 semi hours, our ecr model has a close performance compared to historical real data. Our simulation values are slightly lower after 9:00am, but the differences are within a reasonable range.

In addition to ecr, we also do the same simulation on complete-rate (cr) to model the actual cr which is the boolean variable indicating whether a driver has successfully answered the call and finished the ride. Similar to ecr model, in Figure 3b, our simulation values are bit lower than the real data. But the differences are acceptable when averaged over millions of rides over two weeks.



(a) 8:30am-9:00am, Peak time slot                    (b) 11:30am-12:00pm, Off-peak time slot
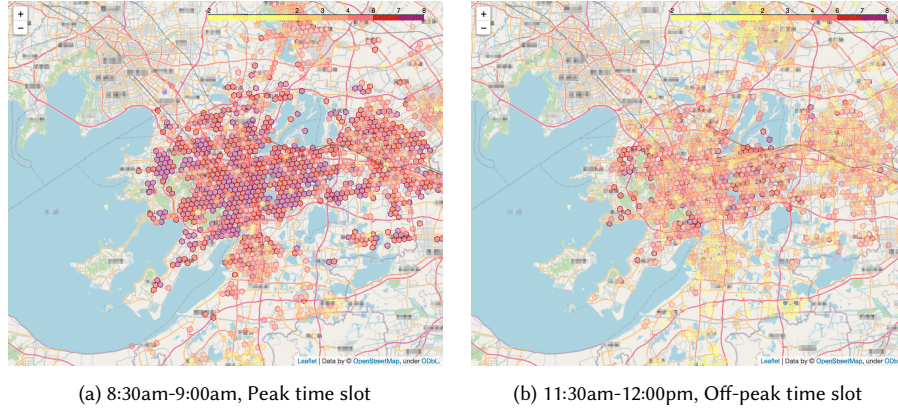
Fig. 5. Spatio-temporal values across the city during peak and off-peak time slots

Table 2.  Sum of destination delta ecr values of grids in short supply

| Algorithms | Historical | IP Policy(ours) | BCQ | BEAR | BCQ+IP | BEAR+IP |
|---|---|---|---|---|---|---|
| $\sum D_t^g$(weekdays) | 302.47 | 625.10 | 307.17 | 299.36 | 671.27 | 506.76 |
| Increase | - | 106.67% | 1.55% | -1.02% | 122% | 67.5% |
| $\sum D_t^g$(weekends) | 72.49 | 182.63 | 102.35 | 102.87 | 200.75 | 139.18 |
| Increase | - | 152% | 41.19% | 41.90% | 177% | 91.99% |

## 5.2   Algorithm Performance

We use one week of historical trip data in one city obtained from a major ride-hailing platform dispatching platform as our training data for learning the Spatio-Temporal network. We build the Spatio-Temporal network with two hidden dense layers and ReLU activation functions. As mentioned in Section 4.6, the Spatio-Temporal network $\mathcal{V}$ is then used to solve an IP policy. When solving the policy, we use $\beta = 0.5$.

To evaluate our policy, we solve the IP policy on a week of historical unseen trip data and compare discounts assigned by our policy, with historical discounts, where historical discounts are taken according to the pricing strategy used at that time. Because our problem is under the batch DRL setting, which is purely off-line with no interactions with the environment, we pick two off-line DRL algorithms BCQ[5] and BEAR[8] for benchmarks. Since BCQ and BEAR do not include an Integer Programming stage to limit the budget for discounting, if we use the same reward definition as mentioned in equation 1, we would result in a policy assigning largest discounts to all trips. Hence, when training BCQ and BEAR algorithms, we add a penalty term to each reward to indicate how much we need to compensate for distributing the discount, the penalty term is balanced by $\alpha$:

$$R_{\hat{Y}} = R_Y - \alpha * (1 - a) * GMV \tag{7}$$

where $a$ represents the discrete discount actions, and takes value of 0.75, 0.8, 0.85, 0.9, 0.95 or 1.



(a) Value of supply minus demand          (b) $D_t^g$ value given by our policy          (c) $D_t^g$ value given by historical data

Fig. 6.  $D_t^g$ comparisons between our policy and historical data (Due to data security issues, we blur out the real grid IDs and names shown on the map.)

Figure 4a shows the distribution over discounts by different policies trained using different algorithms. Historical data has the largest number of rides with no discounts(label = 100), 20% off(label = 80) and 25% off(label = 75). In other words, historical data turns to make more extreme actions to assign larger discounts. Our policy assigns more 5%-off(label = 95) and 10% off(label = 90) discounts to the rides comparing to historical data.On the contrary, our IP policy barely assigns 25% off actions to the trips. Discounts assigned by BCQ are mainly gathered at 5%-off, while BEAR

algorithm has the largest number of 10%-off rides. To study the effect of Integer Programming(IP), we also combine BCQ and BEAR with IP, where we train BCQ and BEAR as normal, and then instead of directly using their policies during evaluation, we train an IP policy using the Q network trained by the algorithms. The results are labeled as BCQ_IP and BEAR_IP. The combination has brought the discount distribution gathering around 5% of BCQ and BEAR down to a level similar to our policy and increase the distribution on other discounts.

In addition to the overall discount distribution, we also analyze the discounts associated with trip fares. In Figure4b and Figure4c, for different trip fare intervals, we plot the discount distribution over the six discrete actions. For example, for rides with fare between 70 and 80, our policy turn to assign 20%-off discounts to more than half of the trips. Meanwhile, historical discounts only assign about 10% of the trips with 20%-off discounts. Also, historical discount distribution assigns least amount of 5% off discounts to short rides with fare between 10 and 20. On the contrary, our policy assigns least amount of 5% off actions to long rides with fare between 70 and 80.

We summarize and compare model performance of our policy, BCQ, BCQ_IP, BEAR and BEAR_IP algorithms to historical data in terms of total revenue in table 1. The results are evaluated over one week of unseen data. Our policy gives 3.7% increase over the historical weekday data and 3.18% increase over weekends. As we expected, other batch DRL algorithms such as BCQ and BEAR, which learn policy networks to make decisions, are basically impossible to precisely control the budget they use. In order to make sure they do not exceed the budget, we tune the penalty weight $\alpha$ in equation 7 so as to discourage the policy from making large discount actions. Intuitively, the performance would be negatively effected. From the experiments, we can see that neither weekends nor weekday performance of BCQ or BEAR is better than historical data. To help with the decision-making process and taking better advantages of the value networks, we take the Q networks learned by BCQ and BEAR, but instead of optimizing the policy network by maximizing expected Q values, we use the Integer Programming method mentioned in section 4 as their behavior policy. The results are labeled as BCQ_IP and BEAR_IP respectively. By taking this modification, we witness improvements over weekdays from -0.96% to 3.7% for BCQ and -0.85% to 3.1% for BEAR and over weekends from -0.65% to 3.17% for BCQ and -0.55% to 2.59% for BEAR.

## 5.3   Visualization on City Maps

In addition to quantitative analysis and results, we also discover interesting phenomenons by visualizing our spatio-temporal value net on the city map. For example, Figure 5 shows the spatio temporal values of each grid across the city, at two different time slots. Figure 5a shows the spatio-temporal value distribution at peak time from 8:30am to 9:00am. Figure5b shows the spatio-temporal values of the same city but at the off-peak hours from 11:30am to 12:00pm. Not surprisingly, during the morning peak time slot, the whole city is more busy and more rides are happening compared to the off-peak time slot. Downtown areas are especially more valuable than other places in the city.

To see how our algorithm helps to balance the supply and demand situation of the city. We introduce a new quantitative feature called Dest_delta_ecr, and denote it as $D$. The feature is calculated for each grid of the city at each time slot t. It is defined as $D_t^g = \sum_{i=1}^{M_t} \delta \text{ecr}_i$, where $M_t$ is the set of all the rides arriving at destination grid $g$ at time slot $t$ and $\delta \text{ecr}_i$ is the change in ecr of each ride $i$ in the set $M_t$ under its assigned discount.

We also calculate supply minus demand information of each hexagon grid of the city map at each time slot t. If supply minus demand is smaller than zero, we can say that at this time slot t, the grid is in short supply and needs more drivers. For example, in Figure 6a, the grid which we are focused on is a residential area. During the morning peak time slot on a weekday, lots of people are leaving the area for work, thus, the grid has a negative value of supply minus

demand, which means it's in short supply. Hence, we want to encourage rides to the area at that time slot, which will also bring drivers to the area to compensate for the short supply.

Under this situation, historical data, as shown in figure 6c has $D_t$ equals 0.246, while our policy in Figure 6b has $D_t$ equals 0.383 at the same grid and same time slot. This means, discount distribution by our policy, compared to historical distribution, assigns more discounts on trips going to that residential grid. In other words, we increase the passengers' ecr of the rides going to the grid at that time slot in order to help with the supply and demand imbalance of the city.

To better analyze and has a broader view of the effectiveness of our policy on the imbalanced supply and demand problem, we sum over all $D_t^g$ only on destination grids that are in short supply. Table 2 summarizes the results. We want the grids in short supply to have more drivers coming to meet the demands, which means we expect the sum over all $D_t^g$ for grids in short supply to be as large as possible but meanwhile do not exceed the budget. For weekdays, our policy has increased the sum of destination delta ecr values by 106.67%, and has increased the sum of destination delta ecr values by 152% for weekends. BCQ and BEAR have also increased about 40% for weekends, but BCQ only does a slightly better job for weekdays. However, when combining with IP, BCQ_IP boosts the performance by 122% for weekdays and 177% for weekends. BEAR_IP also does a much better job than vanilla BEAR. The result indicates the importance of IP and the advantage over a normal policy network in our application of discount distribution.

## 6 CONCLUSION

This paper has proposed an offline deep reinforcement learning method to learn the spatial-temporal long-term value in the whole city, as a guidance to dynamically adjust the spatial-temporal distribution of demand, to improve the efficiency of the marketplace. The goal is to optimize the ride-hailing platform's long-term efficiency, as well as balancing customer demands and driver supply in the city. To do so, we model the pricing problem as a sequential decision-making problem. Experiments demonstrate the effectiveness of the proposed algorithm and the importance of the integer programming step. We deliver an improvements in both total revenue and efficiency in balancing supply and demand across the city.

## REFERENCES

[1] Jie Bao, Tianfu He, Sijie Ruan, Yanhua Li, and Yu Zheng. 2017. Planning Bike Lanes Based on Sharing-Bikes' Trajectories. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Halifax, NS, Canada) *(KDD '17)*. New York, NY, USA, 1377–1386.

[2] Kostas Bimpikis, Ozan Candogan, and Daniela Saban. 2017. Spatial Pricing in Ride-Sharing Networks. In *Proceedings of the 12th Workshop on the Economics of Networks, Systems and Computation (NetEcon '17)*. Association for Computing Machinery, New York, NY, USA, Article 5, 1 pages.

[3] Chuqiao Chen, Fugen Yao, Dong Mo, Jiangtao Zhu, and Xiqun (Michael) Chen. 2021. Spatial-temporal pricing for ride-sourcing platform with reinforcement learning. *Transportation Research Part C: Emerging Technologies* 130 (2021), 103272.

[4] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*. PMLR, 2052–2062.

[5] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*. PMLR, 2052–2062.

[6] Andrey Glaschenko, A. Ivaschenko, G. Rzevski, and P. Skobelev. 2009. Multi-Agent Real Time Scheduling System for Taxi Companies.

[7] Fang He and Zuo-Jun Max Shen. 2015. Modeling taxi services with smartphone-based e-hailing applications. *Transportation Research Part C: Emerging Technologies* 58 (2015), 93–106.

[8] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. 2019. Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction. *Advances in Neural Information Processing Systems* 32 (2019), 11784–11794.

[9] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. *arXiv preprint arXiv:2006.04779* (2020).

[10] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* (2020).

[11] Minne Li, Zhiwei Qin, Yan Jiao, Yaodong Yang, Jun Wang, Chenxi Wang, Guobin Wu, and Jieping Ye. 2019. Efficient ridesharing order dispatching with mean field multi-agent reinforcement learning. In *The World Wide Web Conference*. 983–994.

[12] Qingquan Li, Zhe Zeng, Bisheng Yang, and Tong Zhang. 2009. Hierarchical route planning based on taxi GPS-trajectories. In *2009 17th International Conference on Geoinformatics*. 1–5.

[13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (Feb. 2015), 529–533.

[14] Mehdi Nourinejad and Mohsen Ramezani. 2020. Ride-Sourcing modeling and pricing in non-equilibrium two-sided markets. *Transportation Research Part B: Methodological* 132 (2020), 340–357. https://doi.org/10.1016/j.trb.2019.05.019 23rd International Symposium on Transportation and Traffic Theory (ISTTT 23).

[15] Wenjie Shang, Qingyang Li, Zhiwei Qin, Yang Yu, Yiping Meng, and Jieping Ye. 2021. Partially observable environment estimation with uplift inference for reinforcement learning based recommendation. *Machine Learning* (2021), 1–38.

[16] Wenjie Shang, Yang Yu, Qingyang Li, Zhiwei Qin, Yiping Meng, and Jieping Ye. 2019. Environment reconstruction with hidden confounders for reinforcement learning based recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 566–576.

[17] Xiaocheng Tang, Zhiwei Qin, Fan Zhang, Zhaodong Wang, Zhe Xu, Yintai Ma, Hongtu Zhu, and Jieping Ye. 2021. A Deep Value-network Based Approach for Multi-Driver Order Dispatching. arXiv:2106.04493 [cs.LG]

[18] Yongxin Tong, Yuqiang Chen, Zimu Zhou, Lei Chen, Jie Wang, Qiang Yang, Jieping Ye, and Weifeng Lv. 2017. The Simpler The Better: A Unified Approach to Predicting Original Taxi Demands Based on Large-Scale Online Platforms. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Halifax, NS, Canada) *(KDD '17)*. New York, NY, USA, 1653–1662.

[19] Hado van Hasselt, Arthur Guez, and David Silver. 2015. Deep Reinforcement Learning with Double Q-learning. arXiv:1509.06461 [cs.LG]

[20] Zhaodong Wang, Zhiwei Qin, Xiaocheng Tang, Jieping Ye, and Hongtu Zhu. 2018. Deep Reinforcement Learning with Knowledge Transfer for Online Rides Order Dispatching. In *2018 IEEE International Conference on Data Mining (ICDM)*. 617–626.

[21] Huasen Wu, R Srikant, Xin Liu, and Chong Jiang. 2015. Algorithms with Logarithmic or Sublinear Regret for Constrained Contextual Bandits. *Advances in Neural Information Processing Systems* 28 (2015), 433–441.

[22] Tang Xin-min, Wang Yu-ting, and Han Song-chen. 2010. Aircraft Taxi Route Planning for A-SMGCS Based on Discrete Event Dynamic System modeling. In *2010 Second International Conference on Computer Modeling and Simulation*, Vol. 1. 224–228.

[23] Zhe Xu, Zhixin Li, Qingwen Guan, Dingshui Zhang, Qiang Li, Junxiao Nan, Chunyang Liu, Wei Bian, and Jieping Ye. 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* (London, United Kingdom) *(KDD '18)*. New York, NY, USA, 905–913.

[24] Hai Yang, Chaoyi Shao, Hai Wang, and Jieping Ye. 2020. Integrated reward scheme and surge pricing in a ridesourcing market. *Transportation Research Part B: Methodological* 134 (2020), 126–142.

[25] Huaxiu Yao, Fei Wu, Jintao Ke, Xianfeng Tang, Yitian Jia, Siyu Lu, Pinghua Gong, Jieping Ye, and Zhenhui Li. 2018. Deep multi-view spatial-temporal network for taxi demand prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.

[26] Liteng Zha, Yafeng Yin, and Zhengtian Xu. 2018. Geometric matching and spatial pricing in ride-sourcing markets. *Transportation Research Part C: Emerging Technologies* 92 (2018), 58–75.