

Reinforced Contact Tracing and Epidemic Intervention

Tao Feng, Sirui Song, Tong Xia, Yong Li

Abstract—The recent outbreak of COVID-19 poses a serious threat to people’s lives. Epidemic control strategies have also caused damage to the economy by cutting off humans’ daily commute. In this paper, we develop an Individual-based Reinforcement Learning Epidemic Control Agent (IDRLECA) to search for smart epidemic control strategies that can simultaneously minimize infections and the cost of mobility intervention. IDRLECA first hires an infection probability model to calculate the current infection probability of each individual. Then, the infection probabilities together with individuals’ health status and movement information are fed to a novel GNN to estimate the spread of the virus through human contacts. The estimated risks are used to further support an RL agent to select individual-level epidemic-control actions. The training of IDRLECA is guided by a specially designed reward function considering both the cost of mobility intervention and the effectiveness of epidemic control. Moreover, we design a constraint for control-action selection that eases its difficulty and further improve exploring efficiency. Extensive experimental results demonstrate that IDRLECA can suppress infections at a very low level and retain more than 95% of human mobility.

Index Terms—COVID-19, RL, GNN



1 INTRODUCTION

The recent outbreak of COVID-19 has caused thousands of infections and deaths. Similar to most epidemics that can spread via human contact [1], control the spread of the COVID-19 virus requires cutting off human contacts. Governments have taken different epidemic-control strategies, such as travel-restriction orders, individual quarantine policies, and city lockdown [2]. However, restricting human’s daily mobility and gathering will inevitably pose a negative effect on economic growth. The current epidemic control strategies for COVID-19 has ultimately caused damage to the economy [3], [4].

To control the epidemic both efficiently and effectively, researchers have proposed smart and computational Epidemic-Prevention-and-Control (EPC) strategies in both group level and individual level. Group-level EPC strategies [5], [6] aim to select customized epidemic-control actions for each population group. These works are mainly based on the SIR model [7] which can characterize the development trend of the epidemic from a group-level view. However, Group-level EPC strategies ignore the unique situation of each individual, which may easily cause unnecessary mobility intervention costs or secondary transmission of infection. By contrast, individual-based EPC strategies exploit individual information to estimate infection risk for each individual, and further select a customized epidemic-control action for each individual [8]. However, current individual-based EPC strategies [9], [10], [11], [12], [13] lack a module to estimate the spread of the virus through complex contacts between individuals. To achieve an efficient and effective EPC result, we in this paper aim to maximally make use of

available information and design an individual-based EPC strategy that can both minimize the number of infections and the social cost of epidemic control.

The main challenges of our research are three-fold. *First*, primitive individual information can hardly reflect an individual’s infection risk. For example, an asymptomatic patient who has a very high infection risk is usually hard to detect just through symptoms. In other words, the large population and their complex information form a vast state space for control, making it very hard to extract effective information to support the selection for control actions. *Second*, the large and complex action space exacerbates the difficulty of control-action selection. If there exist M people and d kinds of control actions, the action space is M^d , which is growing exponentially. *Third*, searching for a policy that achieves the dual objective of minimizing both infections and the social cost of implementing the strategy is hard. The two optimization goals will influence each other. For example, better control of the epidemic requires greater control efforts, which will naturally increase mobility intervention costs.

To solve the above challenges, we propose an Individual-based Reinforcement Learning Epidemic Control Agent (IDRLECA) by combining Graph Neural Network and Reinforcement Learning approach. Specifically, to deal with the vast-state-space challenge, we design an infection probability model to calculate the current infection probability of each individual, whose result is further added to the individual’s state as auxiliary information. In order to better extract individual features hidden in his/her daily commute, we design a novel GNN which inputs with individuals’ states their visiting history and estimates their infection risks of individuals. As for the large-action-space challenge, we design and impose a constraint to control-action selection by requiring individuals with larger calculated infection probability should receive more stringent control actions.

- T. Feng, S. Song, T. Xia and Y. Li are with Beijing National Research Center for Information Science and Technology (BNRist), Department of Electronic Engineering, Tsinghua University, Beijing 100084, China. Email: liyong07@tsinghua.edu.cn.

In response to the dual-objective optimization challenge, we carefully design a reward function considering both the social cost of EPC and the effectiveness of infection suppression. More importantly, the reward function is able to efficiently guide training.

We build a simulation environment based on the PAPW Challenge¹, and experimentally compare the performance of expert EPC strategies, winners in the PAPW Challenge, and our proposed IDRLECA. Extensive results show that IDRLECA achieves the best performance for both infection-suppression and mobility-retaining in all three compared scenarios.

In summary, this paper makes the following contributions:

- We propose IDRLECA to minimize the number of infections and the social cost of EPC. IDRLECA achieves the best performance in both infection-suppression and mobility-retaining compared with expert baselines and PAPW winners.
- We propose a method to address the vast-state-space problem in individual-based EPC. Our method includes an infection-probability model and a novel GNN.
- We design and impose a constraint to control-action selection to improve the exploration efficiency of IDRLECA in the extremely large action space.

The remainder of this paper is as follows. We first introduce our problem formulation in Section 2 and introduce our method in Section 3. The experiment results are presented in Section 4. We introduce the related works in Section 5 and conclude our paper in Section 6.

2 PROBLEM FORMULATION AND CHALLENGES

In this section, we formulate the problem of individual-based EPC and discuss the challenges in finding an effective EPC policy.

2.1 Formulation

We consider a within-city epidemic control scenario. The city is assumed to be composed of N areas and has a population of M . Each individual's health status can be: Susceptible, Asymptomatic, Symptomatic, and Recovered. Asymptomatic and Symptomatic individuals are both infected. Each individual will commute between different areas according to some predefined commute rules. When people are staying in the same area, they have a probability to contact each other and their health status will change from Susceptible to Asymptomatic. The Asymptomatic status will Symptomatic after a predefined incubation time. Symptomatic individuals will be sent to the hospital and transit to Recover after a predefined time of treatment. Note that policymakers cannot distinguish between Susceptible individuals and Asymptomatic individuals. The goal of individual-based EPC is to select a control action for each individual in the Susceptible group and Asymptomatic group to minimize the number of infected people and the cost of intervention measures. Specifically, we define four kinds of

control actions: No Intervention, Confine (no contact with people living outside his/her residential area), Quarantine (no stranger contact), Isolate (no contact). The above modeling for epidemic transmission and individual-based EPC actions is shown in Figure 1.

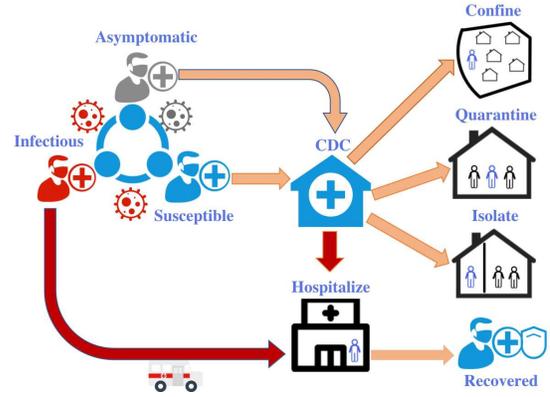


Fig. 1. Epidemic Spread and Intervention(CDC: Center for Disease Control and prevention).

Once the number of infected people exceeds a threshold, the medical system will be penetrated, leading to a rapid increase in medical costs. On the other hand, when the mobility intervention is greater than a certain threshold, the economic system will be paralyzed, also leading to a sharp increase in social cost. So we design the metric *Score* to evaluate the total cost of an EPC policy to consider reducing the infections and maintaining the mobility at the same time. The smaller *Score* value indicates better EPC results. *Score* is defined as follows:

$$Q = \lambda_h * N_h + \lambda_i * N_i + \lambda_q * N_q + \lambda_c * N_c,$$

$$Score = exp \left\{ \frac{I}{\theta_I} \right\} + exp \left\{ \frac{Q}{\theta_Q} \right\},$$

where I denotes the total number of infected people within all simulation days, Q denotes the aggregate of mobility interventions, N_h , N_i , N_q and N_c denote the accumulated number of hospitalized, isolated, quarantined, and confined people for all simulation days, θ_I and θ_Q refers to the soft thresholds for medical system's capacity and economic system's endurance. λ_h , λ_i , λ_q and λ_c denote scale factors.

In this paper, we aim to find an EPC policy that gives daily control actions for all individuals to minimize *Score*.

2.2 Challenges

Finding an effective EPC policy is challenging in three aspects:

2.2.1 Vast State Space

The invisibility of asymptomatic patients and people's complex contacts makes the state space vast. It's difficult to extract effective features for control-action selection. To tackle this challenge, we propose two solutions. We design an infection probability model to calculate the current infection probability of each individual. The probability is added to the state of each individual as auxiliary information. Moreover, IDRLECA employs a novel GNN acquires the

1. PAPW 2020: <https://prescriptive-analytics.github.io/>.

whole individuals' state and the area visited history as input, which can estimate the infection risks through the contact between individuals. The estimated infection risks which measure the individual's ability to potentially infect others are further used as action thresholds to support the selection for actions.

2.2.2 Large Action Space

Individual-level epidemic control aims to select a control action for each individual, which brings an extremely large action space for this control problem. This further leads to low exploration efficiency for reinforcement learning. In order to solve this challenge, we design an infection probability model to calculate the current infection probability for each person and use IDRLECA to output different action control thresholds for each individual. The estimated risks are further used to support RL's selection for action actions.

2.2.3 Dual Objective Optimization

Since our goal is to minimize the social cost *Score* which contains two optimization objectives of the entire epidemic control process. To solve this problem, we propose a special design instant reward, which considers the number of new infections on two consecutive days and the mobility intervention cost on the current day.

3 METHODOLOGY

To tackle the above challenges, we develop an Individual-based Reinforcement Learning Epidemic Control Agent (IDRLECA) that employs a novel GNN and RL approach to search for smart control policies. An overview of IDRLECA is shown in Figure 2. At each time step, IDRLECA collects the health status, intervention state and area-visit-history for each individual and gives each a customized intervention action. In the rest of this section, we will provide the details of the IDRLECA.

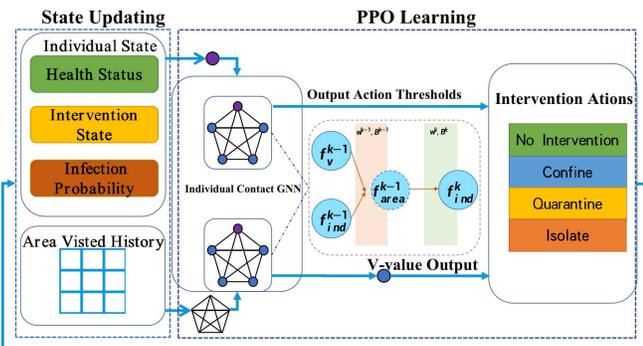


Fig. 2. The detailed structure of proposed IDRLECA.

3.1 Infection Probability Model

The difficulty of epidemic prevention and control lies in how to find asymptomatic infections and how to take timely and effective measures. To help the latter part of IDRLECA efficiently take use of effective information, we here design an infection probability model to estimate the probability of an individual being infected. We define the probability

of infection and health of the i -th person as p_i^{infe} and p_i^{hel} , respectively. The infection probabilities of contacting with strangers and acquaintances are calculated by the simulation environment, denoted as p_s and p_c , respectively. The infection probability model works as follows:

Step 1: Trace back all individuals' area-visit history in the past T time steps.

Step 2: For individual $i, i = 1, 2, \dots, M$, define his/her probability of being healthy as $p_{i,t}^{hel}$ at time step t . $p_{i,0}^{hel}$ is initialized to be 1 if individual i is not infected. we have the following equation to update $p_{i,t}^{hel}$:

$$p_{i,t}^{hel} = p_{i,t-1}^{hel} * (1 - p_s \frac{N_{t-1}^{infe}}{N_{t-1}^{area}}), t = 1, 2, \dots, T, \quad (1)$$

where N_{t-1}^{infe} and N_{area} refer to the number of discovered infections and total number of visitors to the same area as individual i , respectively.

Step 3: Update $p_{i,T}^{hel}$ for acquaintances' contacts:

$$\hat{p}_{i,T}^{hel} = p_{i,T}^{hel} * (1 - p_c). \quad (2)$$

Step 4: Acquire infection probability:

$$p_i^{infe} = 1 - \hat{p}_{i,T}^{hel}, \quad (3)$$

After the above steps, we can obtain the estimated probability of an individual being infected. We will use it as auxiliary information and add it to each individual's state. Also, the estimated probabilities are used as prior knowledge for the agent selection control actions. Note that this estimation is not 100% accurate since our estimation simplifies the process of contact and spread of virus between people. In the later part, we will combine GNN to solve this problem.

3.2 Reinforcement Learning

We propose IDRLECA to search smart strategies to minimize the spread of epidemic and cost of intervention at the same time. We treat all individuals in the area as one agent. Therefore, for IDRLECA, its status and actions are for all people. We use one day as the decision time interval. In the following, we will introduce our design of state, action and reward:

- **State:** The state of IDRLECA is the integration of each individual's information, which is obtained at the start of one day. For each individual, the state includes infection state, intervention state, and the probability of infection calculated by Equation (1)~(3).
- **Action:** The action at each step for the agent is to determine the intervention measure of each individual. The action contains no intervention, confine, and quarantine. In order to ensure the flexibility of the policy, we set the implementation time of actions to one day.
- **Reward:** The goal of our method is to minimize the total number of infected people, and to minimize the total intervention cost at the same time. Considering

our dual objective optimization, we set the reward r as follows,

$$r = -\exp\left\{\frac{\Delta I}{\theta_I}\right\} - \exp\left\{\frac{\Delta Q}{\theta_Q}\right\}, \quad (4)$$

where ΔI and ΔQ denote the daily incremental part of the number in the infected population between consecutive days and the cost of mobility intervention on the day, respectively.

- **Learning Algorithm:** IDRLECA employs a Proximal Policy Optimization (PPO) [14] agent to find the optimal strategy that minimizes the number of infected people and the cost of prevention at the same time. The PPO agent adopts the actor-critic framework. The critic network is used to estimate the long-term reward of the action, and the actor network is to find the optimal action policy to achieve dual objective optimization. We also add an entropy bonus to ensure sufficient exploration when RL training [14].

3.3 Individual Contact GNN

Since asymptomatic patients are indistinguishable, it's hard to trace all the contacts and infections caused by them. Moreover, vast modern traffic and complex social network structure make it more challenging to estimate the infection risk of each individual. To deal with this challenge, we propose a novel GNN, namely Individual Contact GNN, to estimate the infection risk of each individual. Individual Contact GNN is used to build both the actor network and critic network in IDRLECA. The GNN regards individuals and city areas as two kinds of nodes. This enables us to model individual-individual contacts by individual-area-individual contacts, which further helps us to avoid the extremely large individual-individual contact matrix (size $M * M$).

Specifically, Individual Contact GNN is designed on the basis of GraphSage [15]. The state input to the GNN consists of health status, intervention state, infection probability for all individuals, and the edge-information inputs are the area-visit-history at different time steps. We use f_{area}^k, f_{ind}^k to denote for the area-nodes' features and individual-nodes' features outputted by the k -th GNN layer, respectively. The detailed layer-calculation of Individual Contact GNN is as follows:

$$f_c^{k-1} = \text{softmax}(f_v^{k-1}), \quad (5)$$

$$f_{area}^{k-1} = \sigma(W^{k-1}(f_c^{k-1})^T f_{ind}^{k-1} + B^{k-1}), \quad (6)$$

$$f_{ind}^k = \sigma(W^k f_c^{k-1} f_{area}^{k-1} + B^k), \quad (7)$$

where f_v^{k-1} denotes for the area's visit history at the $k-1$ time step, $W^{k-1}, B^{k-1}, W^k, B^k$ denotes for trainable parameters.

In the above equations, Equation (5) uses the area-visit-history as edge weights; Equation (6) aggregates weighted visitors' characteristics to calculate the area-node feature; Equation (7) aggregates the features of areas where an individual has visited to calculate individual-node feature.

3.4 Constraint for Control-Action Selection

As discussed before, the EPC problem has a extremely large action space, which challenges policy search. To address this issue, we incorporate prior knowledge into the control-selection step. Specifically, we let the actor network of IDRLECA first outputs four values $\langle p_{i,1}, p_{i,2}, p_{i,3}, p_{i,4} \rangle$ for individual $i, i = 1, 2, 3, \dots, M$. Then, we transform the four values to three thresholds:

$$P_{i,1} = \frac{e^{-p_{i,1}}}{e^{-p_{i,1}} + e^{-p_{i,2}} + e^{-p_{i,3}} + e^{-p_{i,4}}}, \quad (8)$$

$$P_{i,3} = \frac{e^{-p_{i,1}} + e^{-p_{i,2}}}{e^{-p_{i,1}} + e^{-p_{i,2}} + e^{-p_{i,3}} + e^{-p_{i,4}}}, \quad (9)$$

$$P_{i,2} = \frac{e^{-p_{i,1}} + e^{-p_{i,2}} + e^{-p_{i,3}}}{e^{-p_{i,1}} + e^{-p_{i,2}} + e^{-p_{i,3}} + e^{-p_{i,4}}}. \quad (10)$$

Through the above equations, we can ensure $0 \leq P_{i,1} \leq P_{i,2} \leq P_{i,3} \leq 1$. Thus, $P_{i,1}, P_{i,2}, P_{i,3}$ can be used as different infection risk levels, which considers the risk of individual infection and individual's ability to potentially infect others. It's natural and reasonable to expect that an individual with a higher infection risk should receive a more stringent control action. The infection risk levels are further used as the thresholds for the infection probabilities estimated in Section 2, which imposes a constraint that individuals with higher infection probability will have higher infection risk and receive more stringent control actions. In this way, individuals with high probability of infection are not identified as low risk, thus reducing unnecessary strategy exploration.

By comparing the pre-calculated infection probability p_i^{inf} with $\langle P_{i,1}, P_{i,2}, P_{i,3} \rangle$, we define the action-selection rule in Table 1. It can be seen from Table 1 that as the infection probability goes from low to high, the corresponding intervention actions become more and more stringent. There are different thresholds for different individuals, which fully takes into account the differences in individual states.

TABLE 1
Action-Selection Rule.

Infection probability	Intervention actions
$0 \leq p_i^{inf} \leq P_{i,1}$	No intervention
$P_{i,1} \leq p_i^{inf} \leq P_{i,2}$	Confine
$P_{i,2} \leq p_i^{inf} \leq P_{i,3}$	Quarantine
$P_{i,3} \leq p_i^{inf} \leq 1$	Isolate

3.5 Avoiding extreme experiences

Similar to DURLECA where RL is used for epidemic control [6], it is possible to encounter extreme states or actions during the RL exploration in IDRLECA's training. This may severely impact exploration efficiency and result in local optimums. Inspired by DURLECA, we have a rule to avoid these extreme experiences:

- The infection-increase threshold I_t : During the agent's exploration process, if the number of new infections on a certain day exceeds I_t , the current episode will be stopped and a large penalty will be given to the reward of the agent.

4 EXPERIMENTS

In this section, we conduct extensive experiments on four scenarios to answer the following research questions:

- **RQ1:** Can IDRLECA minimize the number of infections and the cost of interventions?
- **RQ2:** Can IDRLECA be adapted to different scenarios?
- **RQ3:** How does IDRLECA compare to expert policies and PAPW winners?

4.1 Experiment Setup

In the following, we introduce more details about our experiment design.

4.1.1 Simulation Environment

We build a simulation environment mainly based on the PAPW Challenge². The simulated disease has an $R0$ range from 2 to 2.5, which is similar to COVID-19³. The total simulation time is 60 days. Every individual has a pre-defined commute pattern. To simulate a more practical EPC scenario, we add a new rule in the original simulator: all symptomatic patients should be sent to the hospital.

4.1.2 Comparison Scenarios

We define t_{start} as the days to start epidemic intervention after discovering the first patient.

- **Scenario-Default:** $N = 11, M = 10000, t_{start} = 1$. This scenario is to verify the EPC performance of IDRLECA in an ordinary epidemic scenario.
- **Scenario-Larger:** $N = 98, M = 10000, t_{start} = 1$. This scenario is to verify whether IDRLECA is suitable for scenarios with greater individual mobility.
- **Scenario-Changeable:** $N = 11, M = 10000, t_{start} = 1$. Compared with Scenario-Default, people's commute patterns are more changeable in this scenario. This scenario is to verify whether IDRLECA is applicable when there are greater differences in individuals' characteristics.
- **Scenario-Late:** $N = 11, M = 10000, t_{start} = 5$. Compared with Scenario-Default, this scenario starts intervention after 5 days of discovering the first patient. This scenario is to verify the EPC performance of IDRLECA with a late intervention.

4.1.3 Evaluation Metrics

- I : The total number of infected people in all simulation days. It is used to measure the effectiveness of EPC strategies in suppressing infections.
- Q : The aggregated mobility interventions defined in Section 2. To have a fair comparison with PAPW winners, we set $\lambda_h = 1, \lambda_i = 0.5, \lambda_q = 0.3$ and

$\lambda_c = 0.2$, which are the same with the setting in the PAPW Challenge.

- **Score:** The social cost of epidemic control policy which is defined in Section 2. We set $\theta_I = 500$ and $\theta_Q = 10000$, which are the same with the setting in the PAPW Challenge.

4.1.4 Comparison Baselines

We set up 4 expert baselines to simulate EPC strategies in the real world:

- *No Intervention:* No intervention at all.
- *Lockdown* [2]: Lockdown the city for successive 60 days.
- *Expert(0.01)* and *Expert(0.015)*: Baselines based on the infection probability model. Isolate individuals whose infection probability is higher than a given threshold.

We compare DIRLECA with two baselines commonly used in epidemic research:

- *Degree - Sample* [13]: If the number of an individual's acquaintances n is more than 4, isolate the individual with a probability $(n - 4)/n$.
- *Degree - Order* [12]: Count the number of contacts of an individual in the past 5 days. Select the top 30% for isolation.

We compare IDRLECA with PAPW winners:

- *GBM* [9]: a baseline for epidemic intervention by predicting individual health states, which strikes a balance between precision and recall.
- *EITL* [10]: a heuristic baseline that adjusts the epidemic strategy through a heuristic algorithm, which based on evaluating the intervention action effectiveness and understanding resulting patterns and interpret causality.
- *HRLI* [11]: a state-of-the-art RL baseline combining individual prevention with regional control.

4.2 Results Analysis

We compare IDRLECA with all the baselines when $t_{start} = 1day$. Table 2 shows our main results. IDRLECA is better than all baselines in three scenarios in metric Score. For instance, compared with the best baseline HRLI in Scenario-Default, our method can reduce the number of infected persons by 26.73% and the cost of mobile intervention by 34.12%.

In the four expert baselines, we can find that *No Intervention* will aggravate the spread of infectious diseases and eventually lead to the paralysis of the medical system. The other three expert baselines can limit the spread of epidemic to some extent. However, these strategies have paid huge mobile intervention costs in order to reduce the number of infections, thereby greatly increasing the total social cost. For example, *Lockdown* is a common method in our real life when dealing with epidemic, which achieves best performance in minimizing infections at the expense of the maximum mobile intervention cost. Compared with the four expert baselines, IDRLECA can minimize the

2. PAPW 2020: <https://prescriptive-analytics.github.io/>. Simulator: <https://hzw77-demo.readthedocs.io/en/round2/>.

3. World Health Organization. (2020, May 8). Report of the WHO-China Joint Mission on Coronavirus Disease 2019 (COVID-19). Retrieved May 8, 2020, from: <https://www.who.int/docs/default-source/coronaviruse/who-china-joint-mission-on-covid-19-final-report.pdf>

TABLE 2
Performance comparison in three scenarios when $t_{start} = 1days$

Method	Scenario-Default			Scenario-Larger			Scenario-Changeable		
	I	Q	Score	I	Q	Score	I	Q	Score
No Intervention	8289	123153.00	>10000	6588	92563.00	>10000	8115	113596.00	>10000
Lockdown [2]	58	294460.50	>10000	56	294508.50	>10000	55	294491.50	>10000
Expert(0.01)	276	6997.50	3.75	204	9187.50	4.01	294	7837.00	3.99
Expert(0.015)	319	8210.00	4.16	269	8404.50	4.03	328	8724.50	4.32
Degree-Sample [13]	1108	212940.00	>10000	1212	211146.50	>10000	943	212498.00	>10000
Degree-Order [12]	3557	120731.00	>10000	2302	92569.50	>10000	3133	119958.50	>10000
GBM [9]	210	6408.21	3.42	177	4794.87	3.04	193	6091.76	3.31
EITL [10]	220	7067.01	3.58	190	5640.15	3.22	205	7899.03	3.71
HRLI [11]	187	5689.79	3.22	183	4935.03	3.08	187	7112.14	3.49
IDRLECA	137	3748.58	2.77	170	4606.17	2.99	153	4068.09	2.86

TABLE 3
Performance comparison in Scenario-Late when $t_{start} = 5days$

Method	Scenario-Late		
	I	Q	Score
No Intervention	8040	119175.00	>10000
Lockdown [2]	70	274364.50	>10000
Expert(0.01)	340	8985.50	4.43
Expert(0.015)	323	8388.00	4.22
Degree-Sample [13]	2091	195949.00	>10000
Degree-Order [12]	3331	115858.00	>10000
GBM [9]	304	7808.13	4.02
EITL [10]	291	8193.50	4.06
HRLI [11]	270	7197.86	3.77
IDRLECA	193	5061.64	3.13

infections and retain large amounts of mobility at the same time.

Compared to *GBM*, *EITL* and two baselines commonly used in epidemic research, *IDRLECA* performs better than them mainly because it considers more individual characteristics and the long-term impact of current actions when making decisions.

Compared to *HRLI*, we find that *IDRLECA* outperforms them in all metrics which may be because the GNN in our method models the contact between individuals and estimates individual infection risks through contact. *IDRLECA* can find hidden infections through GNN and thus be able to stop the spread of epidemic quickly at minimal mobility intervention cost, which will be verified in Case Study.

From the results of Scenario-Larger and Scenario-Changeable, we can find *IDRLECA* can still guarantee the minimum number of infections and mobility intervention costs in more changeable and flexible scenarios.

Late intervention to an epidemic is very common in the real world. An effective control strategy should be able to stop the spread of the epidemic in time with the least cost of mobility intervention in the case of late intervention. We perform our experiment in Scenario-Late, and the results are shown in Table 3. The results show that our method performs best in metric Score compared with other baselines in the case of late intervention.

In Figure 3, we compared the number of infections and the cost of mobility intervention between *IDRLECA* and the best baseline method *HRLI* in Scenario-Late with $t_{start} = 5days$ within 60 days. It can be found that our method can not only stop the spread of epidemic diseases faster, but also reduce the cost of intervention during the peak period of the epidemic.

In order to verify the effectiveness of our method for individual epidemic prevention and control, we randomly select 100 individuals in Scenario-Default, and draw a heat map of the infection probability change within 60 days in Figure 4. It can be found that the infection probability of 100 people reaches its peak in about 15 days, but soon under the influence of intervention measures by *IDRLECA*, the probability of infection is soon reduced to 0 around the 40th day.

4.3 Case Study

In order to verify the effectiveness of our method in individual prevention and control, we conduct two case studies.

Evaluating individual intervention: To verify the specific effects of our method on individual intervention, we draw the infection probability of a person and the changes in prevention and control measures within 60 days of Scenario-Default in Figure 5. It can be found that our method is very sensitive to the action control of different infection probabilities and can effectively reduce the risk of infection.

Finding hidden infections: In order to verify whether our method can discover hidden infections, we used *IDRLECA* to output actions to individuals with ID 927 and 959 on the 20th day of Scenario-Default: quarantine and confine. However, the infection probability of these two individuals is 0.004 and 0.34 respectively, whose numerical order is exactly opposite to the prevention and control level. We further found that the first person had more contacts and acquaintances in the past five days than the second person. This is because the infection probability calculated in Sec 3.1 only considers the impact of the current discovered infections and simplifies the spread of the epidemic by individual contacts. Our GNN models the contact between individuals and can estimate the individual's potential risk of infection and the ability to potentially infect others. Therefore, although the first person is relatively low in the probability of infection, our model takes into account

the infection risk which measures the harm and risk of secondary transmission of potentially infected individuals, so more strict measures are taken for the first person. Two days later, the first person was detected as infected during the intervention period, which also verifies our findings.

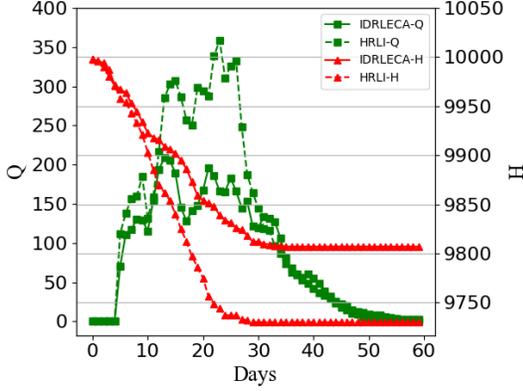


Fig. 3. Q (the aggregated mobility interventions defined in Section 2) and H (the number of healthy people) are changing over time.

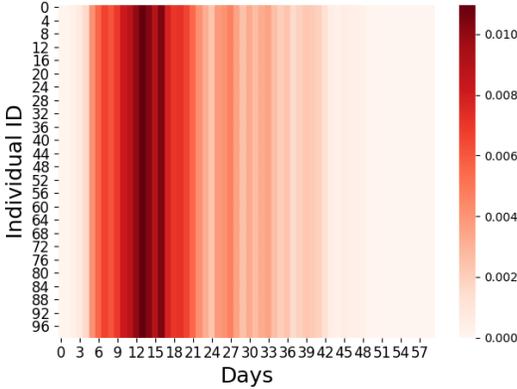


Fig. 4. Change in infection probability of 100 individuals within 60 days(Scenario-Default).

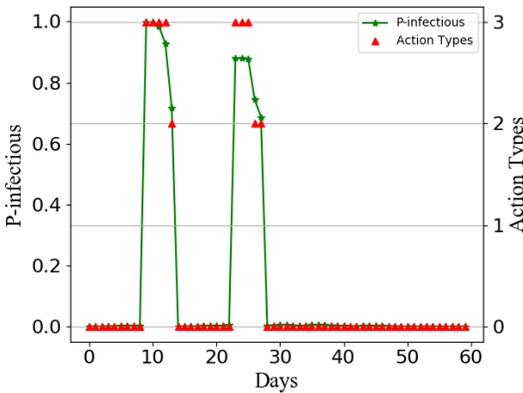


Fig. 5. The relationship between infection probability and intervene action types.

4.4 Ablation Study

To evaluate the effectiveness of our proposed Individual Contact GNN and RL exploration strategy(Avoiding extreme experiences), we take ablation study in this section. We select three baselines and perform experiments in two scenes. No Intervene is the baseline of the blank control. RL-NoGraph and RL-NoEP denote removing GNN and RL exploration strategy(Avoiding extreme experiences) compared with IDRLECA. The results in Table 4 show that removing the GNN network structure will make it difficult for RL to find hidden infections, which will increase the number of infections and the cost of prevention and control. The removal of the exploration strategy(Avoiding extreme experiences) will make it hard for RL to further reduce the number of infections and the cost, falling into a local optimum. Compared with RL-NoGraph and RL-NoEP, IDRLECA can better find hidden infections with the help of GNN, and ensure reasonable and effective exploration under the exploration strategy, so that it can learn better results.

TABLE 4
Ablation study

Method	Scenario-Default			Scenario-Default($t_{start} = 5days$)		
	I	Q	Score	I	Q	Score
No Intervene	8289	123153.00	>10000	8040	119175.00	>10000
RL-NoGraph	200	6606.32	3.42	313	7699.32	4.03
RL-NoEP	192	5779.26	3.25	285	6678.26	3.82
IDRLECA	137	3748.58	2.77	193	5061.64	3.13

5 RELATED WORKS

5.1 Individual-based Infection Simulation and Control Model:

Individual-based Infectious Diseases Model(IBIDM) is an epidemiology model that has emerged in recent years [16]. Compared with traditional infectious disease models, IBIDM can reflect the heterogeneity of individuals and reflect individual-level behavior dynamics, thus more precisely reflect the spread of the epidemic. IBIDM models each individual as a unit, and measures the contact relationship between individuals through social contact network. The Los Alamos Laboratory in the United States has developed an individual-based infectious disease simulation tool, called EpiSimS system, which can effectively simulate the spread, prevention and control of the epidemic based on individual characteristics [17]. Later, some researchers propose Epifast to simulate the spread of Ebola in West Africa, which has higher prediction accuracy and simulation preciseness than traditional methods [18]. There are also many researches related to epidemic control based on these epidemic simulation. [19] studies the trade off between spread of COVID-19 and economic impact and proposes some mechanisms based on group scheduling to strike a balance between epidemic control and economic development. [12], [13], [20] regard individuals as nodes of the graph, and the connections between individuals as edges, and find the individuals who need to be isolated through graphs. [21] introduces mean-field models and complex networks to solve the individual prevention and control of the epidemic.

However, current researches are hard to to effectively extract the status of individuals and their strategies are often unable to cope with various scenarios and conditions. They often only pay attention to the current prevention and control effects, and do not care about the long-term impact of the current decision-making. Therefore, we develop IDRLECA which considers not only how to track and control the infectious and asymptomatic based on the infection status of individuals, but also how to achieve better epidemic prevention while minimizing economic losses in the long term.

5.2 Graph Neural Network for Individual Contacts:

Graph Neural Networks (GNN) are mainly used for node prediction, link prediction and graph prediction tasks. Node prediction refers to predicting the type of a given node [22], [23]. Link prediction means predicting the connection status of two given nodes [24], [25]. Graph prediction aggregates all node features in the graph as the graph feature, and then classifies the type of the graph based on it [26], [27]. There are some commonly used GNN methods. GCN uses the adjacency matrix of nodes as input to learn the relationship between nodes [22]. GAT introduces an attention mechanism on the basis of GNN [28]. GraphSage learns node relationships by aggregating information from neighbor nodes [15].

However, current GNN methods lack a framework to model the spread of epidemic between individuals over a dynamic graph. Therefore, we propose a novel GNN structure to characterize the epidemic-spreading between individuals, whose nodes and edges represent the state features of individuals and contacts between individuals respectively.

5.3 PPO:

PPO algorithm is a new type of policy gradient algorithm and has been applied in many aspects. [14] proposes that PPO strikes a balance between implementation simplicity, sample complexity and difficulty of tuning and achieves good results in many games. [29], [30] proves that PPO can perform well in solving some problems with large-scale and complex state and action space.

Since the PPO algorithm has good stability and adaptability, and can achieve good results in large-scale state and action space problems, we choose PPO as RL algorithm in our problem.

6 CONCLUSION

In this paper, we propose IBRLECA that employs a novel GNN and RL approach to minimize infections as well as the mobility intervention cost in EPC. The proposed GNN can estimate the spread of the virus through contacts between individuals. The training of IBRLECA is guided by a specially designed reward. We design and impose a constraint for control-action selection that eases its difficulty and further improve exploration efficiency. Extensive experiments are conducted on different scenarios to show the effectiveness of our proposed method.

7 APPENDIX

To help reproduce the results, here we present the details of the simulator⁴ and experiment settings.

7.1 Introduction of Simulator

The simulator contains a human mobility model and a disease transmission model. The simulator uses these two models to simulate individuals' movements and the spread of the epidemic among individuals. The two models are briefly introduced below:

Human Mobility Model: The human mobility model simulates individual mobility in a city of N areas with M people. Each area is assumed to belong to one of the three categories: working, residential, and commercial. An individual is associated with two fixed areas: a residential area and a working area. We assume that an individual has different modes of mobility during weekdays and weekends. On weekdays, an individual will move from his/her residential area to his/her working area. After work, he/she may visit a nearby commercial area and then will return to his/her residential area. On weekends, an individual will visit a random commercial area. After that, he/she will return to the residential areas.

Disease Transmission Model: The disease can transmit from an infected individual through acquaintance contacts and stranger contacts. Contacts happen among people within the same region. The infection probabilities of contact with acquaintances and strangers are P_c and P_s , respectively. The disease transmission is simulated every hour.

7.2 Experiment Setting

We set the infection probabilities of contacting with strangers $p_s = 0.01$ and infection probabilities of contacting with acquaintances $p_c = 0.05$. The estimated R_0 is 2-2.5. For the extreme-experience policy, we set $Q_t = 250$. For $Score$, we set $\lambda_h = 1$, $\lambda_i = 0.5$, $\lambda_q = 0.3$ and $\lambda_c = 0.2$, which are the same with the setting in the PAPW Challenge. For the reward and $Score$, we set $\theta_I = 500$ and $\theta_Q = 10000$.

In the training process, the beginning state of an episode is random every time. We train IDRLECA for 200,000 steps, using Adam optimizer with learning rate 0.0001. During testing, the initial setting is fixed in both IDRLECA and the baseline methods. Taking into account the randomness of the simulator, we compared the average results of all methods tested with three random seeds.

7.3 Privacy and implementation issues

Privacy issues: Each area's visited history, the total number of people visiting a particular area and person-person relationship [31] can be obtained by individuals' trajectories. In practical system, user anonymity can be used to reduce the risk of privacy leakage for the individual trajectories and the health history data.

Implementation issues: Our system runs on a central server instead of individuals' smartphones and usually the policymaker has the ability to collect the data needed in our

4. PAPW 2020: <https://prescriptive-analytics.github.io/>. Simulator: <https://hzw77-demo.readthedocs.io/en/round2/>.

model. Besides, some recent techniques like Apple and Google APIs⁵ can be used to collect data without using private information. In the practical implementation of our method, distributed servers and federated learning can be used to protect privacy. The city will be divided into small areas. In each area we have a distributed server that receives encrypted data from smartphones and conducts federated learning with the central server. After training, each distributed server pulls the model from the central server and send reminders to users' smartphones.

ACKNOWLEDGMENTS

This work was supported in part by The National Key Research and Development Program of China under grant 2018YFB1800804, the National Nature Science Foundation of China under U1936217, 61971267, 61972223, 61941117, 61861136003, Beijing Natural Science Foundation under L182038, Beijing National Research Center for Information Science and Technology under 20031887521, and research fund of Tsinghua University - Tencent Joint Laboratory for Internet Innovation Technology.

REFERENCES

- [1] D. Balcan, B. Gonçalves, H. Hu, J. J. Ramasco, V. Colizza, and A. Vespignani, "Modeling the spatial spread of infectious diseases: The global epidemic and mobility computational model," *Journal of computational science*, vol. 1, no. 3, pp. 132–145, 2010.
- [2] T. Hale, A. Petherick, T. Phillips, and S. Webster, "Variation in government responses to covid-19," *Blavatnik school of government working paper*, vol. 31, 2020.
- [3] G. Bonaccorsi, F. Pierri, M. Cinelli, A. Flori, A. Galeazzi, F. Porcelli, A. L. Schmidt, C. M. Valensise, A. Scala, W. Quattrociocchi, and F. Pammolli, "Economic and social consequences of human mobility restrictions under covid-19," *Proceedings of the National Academy of Sciences*, vol. 117, no. 27, pp. 15 530–15 535, 2020. [Online]. Available: <https://www.pnas.org/content/117/27/15530>
- [4] S. Barua *et al.*, "Understanding coronanomics: The economic implications of the coronavirus (covid-19) pandemic," *SSRN Electronic Journal* <https://doi.org/10/ggq92n>, 2020.
- [5] Z. Yang, Z. Zeng, K. Wang, S.-S. Wong, W. Liang, M. Zanin, P. Liu, X. Cao, Z. Gao, Z. Mai *et al.*, "Modified seir and ai prediction of the epidemics trend of covid-19 in china under public health interventions," *Journal of Thoracic Disease*, vol. 12, no. 3, p. 165, 2020.
- [6] S. Song, Z. Zong, Y. Li, X. Liu, and Y. Yu, "Reinforced epidemic control: Saving both lives and economy," 2020.
- [7] W. O. Kermack and A. G. McKendrick, "A contribution to the mathematical theory of epidemics," *Proceedings of the royal society of london. Series A, Containing papers of a mathematical and physical character*, vol. 115, no. 772, pp. 700–721, 1927.
- [8] L. E. Rocha and N. Masuda, "Individual-based approach to epidemic processes on arbitrary dynamic contact networks," *Scientific reports*, vol. 6, p. 31456, 2016.
- [9] S. G. Rizzo, "Balancing precision and recall for cost-effective epidemic containment," [EB/OL], 2020, <https://prescriptive-analytics.github.io/file/3-strizzo.pdf>.
- [10] J.-S. Kim, H. Jin, and A. Züfle, "Expert-in-the-loop prescriptive analytics using mobility intervention for epidemics," 2020.
- [11] Y. Dong, C. Yu, and L. Xia, "Hierarchical reinforcement learning for epidemics intervention," 2020.
- [12] S. Eubank, H. Guclu, V. A. Kumar, M. V. Marathe, A. Srinivasan, Z. Toroczkai, and N. Wang, "Modelling disease outbreaks in realistic urban social networks," *Nature*, vol. 429, no. 6988, pp. 180–184, 2004.
- [13] D. J. Watts and S. H. Strogatz, "Collective dynamics of 'small-world' networks," *nature*, vol. 393, no. 6684, pp. 440–442, 1998.
- [14] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [15] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in neural information processing systems*, 2017, pp. 1024–1034.
- [16] G. J. Milne, J. K. Kelso, H. A. Kelly, S. T. Huband, and J. McVernon, "A small community model for the transmission of infectious diseases: comparison of school closure as an intervention in individual-based models of an influenza pandemic," *PloS one*, vol. 3, no. 12, p. e4005, 2008.
- [17] S. M. Mniszewski, S. Y. Del Valle, P. D. Stroud, J. M. Riese, and S. J. Sydorik, "Episims simulation of a multi-component strategy for pandemic influenza," in *Proceedings of the 2008 Spring simulation multiconference*, 2008, pp. 556–563.
- [18] K. R. Bisset, J. Chen, X. Feng, V. A. Kumar, and M. V. Marathe, "Epifast: a fast algorithm for large scale realistic epidemic simulations on distributed memory systems," in *Proceedings of the 23rd international conference on Supercomputing*, 2009, pp. 430–439.
- [19] J. Augustine, K. Hourani, A. R. Molla, G. Pandurangan, and A. Pasic, "Economy versus disease spread: Reopening mechanisms for covid 19," *arXiv preprint arXiv:2009.08872*, 2020.
- [20] P. S. Park, J. E. Blumentstock, and M. W. Macy, "The strength of long-range ties in population-scale social networks," *Science*, vol. 362, no. 6421, pp. 1410–1413, 2018.
- [21] Q. Wu and T. Hadzibeganovic, "An individual-based modeling framework for infectious disease spreading in clustered complex networks," *Applied Mathematical Modelling*, 2020.
- [22] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [23] Z. Liu, C. Chen, X. Yang, J. Zhou, X. Li, and L. Song, "Heterogeneous graph neural networks for malicious account detection," in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 2018, pp. 2077–2085.
- [24] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," in *Advances in Neural Information Processing Systems*, 2018, pp. 5165–5175.
- [25] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, "Graph convolutional neural networks for web-scale recommender systems," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 974–983.
- [26] D. Bacciu, F. Errica, and A. Micheli, "Contextual graph markov model: A deep and generative approach to graph processing," *arXiv preprint arXiv:1805.10636*, 2018.
- [27] X. Geng, Y. Li, L. Wang, L. Zhang, Q. Yang, J. Ye, and Y. Liu, "Spatiotemporal multi-graph convolution network for ride-hailing demand forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 3656–3663.
- [28] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [29] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Denison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse *et al.*, "Dota 2 with large scale deep reinforcement learning," *arXiv preprint arXiv:1912.06680*, 2019.
- [30] D. Ye, Z. Liu, M. Sun, B. Shi, P. Zhao, H. Wu, H. Yu, S. Yang, X. Wu, Q. Guo *et al.*, "Mastering complex control in moba games with deep reinforcement learning," in *AAAI*, 2020, pp. 6672–6679.
- [31] K. Xu, K. Zou, Y. Huang, X. Yu, and X. Zhang, "Mining community and inferring friendship in mobile social networks," *Neurocomputing*, vol. 174, pp. 605–616, 2016.