

# GAC: A Deep Reinforcement Learning Model Toward User Incentivization in Unknown Social Networks

Shiqing Wu<sup>a,\*</sup>, Weihua Li<sup>b</sup>, Quan Bai<sup>a</sup>

<sup>a</sup>University of Tasmania, Australia

<sup>b</sup>Auckland University of Technology, New Zealand

---

## Abstract

In recent years, providing incentives to human users for attracting their attention and engagement has been widely adopted in many applications. To effectively incentivize users, most incentive mechanisms determine incentive values based on users' individual attributes, such as preferences. These approaches could be ineffective when such information is unavailable. Meanwhile, due to the budget limitation, the number of users who can be incentivized is also restricted. In this light, we intend to utilize social influence among users to maximize the incentivization. By directly incentivizing influential users in the social network, their followers and friends could be indirectly incentivized with fewer incentives or no incentive. However, it is difficult to identify influential users beforehand in the social network, as the influence strength between each pair of users is typically unknown. In this work, we propose an end-to-end reinforcement learning-based framework, named Geometric Actor-Critic (GAC), to discover effective incentive allocation policies under limited budgets. More specifically, the proposed approach can extract information from a high-level network representation for learning effective incentive allocation policies. The proposed GAC only requires the topology of the social network and does not rely on any prior information about users' attributes. We use three real-world social network datasets to evaluate the performance of the proposed GAC. The experimental results demonstrate the effectiveness of the proposed approach.

*Keywords:* Incentive Allocation, Reinforcement Learning, Unknown Social Network, Social Simulation

---

## 1. Introduction

In recent years, incentive mechanisms used to incentivize users to select behaviors that are beneficial to the incentive providers have been widely studied and applied in many fields. Typical examples can be promoting sales [1, 2], hiring workers [3, 4, 5], encouraging beneficial behaviors [6, 7]. Such processes are computationally

---

\*Corresponding Author. Email: shiqing.wu@utas.edu.au

modeled as *incentive allocation* problem, where the goal is to incentivize users with effective incentives under a budget limitation, such that the number of users who take the behavior that the incentive provider expects is maximized.

In general, an incentive provider tends to minimize the cost of providing an incentive such that the utility derived from the incentive can be maximized. Whereas users may not accept the provided incentives, since they expect to maximize their individual profit as well. This conflict existing between the incentive provider and users makes incentive allocation a challenging problem, as overpricing the incentive would waste the budget while underpricing could fail to incentivize the user [8]. Hence, one key factor of successful incentive allocation is whether the incentive structure and the pricing policies for users are reasonable or not. More specifically, the pricing policies should consider the utilities of both incentive providers and users. Some studies attempted to model users' attributes, such as users' preferences and skill abilities [4, 7, 9], to generate "optimal" incentives. However, these methods could be impotent when such information is lacking or unavailable. In this case, it is necessary to propose an effective approach that does not rely on prior knowledge about users to distribute incentives.

Meanwhile, the limited budget implies that not all users could receive sufficient incentives, i.e., the overall result of user incentivization would be restricted. Some studies also consider exploiting social influence to facilitate user incentivization, as information diffusion plays a crucial role in propagating persuasive information among friends in a social network [10]. For example, many business companies have advertised their products and brand by providing free products to a few influential users on social media, hoping they can promote products to their friends and followers [11]. Users' decision-making can be influenced by their neighbors or other influential users in the social network [12]. In other words, if we can appropriately incentivize some influential users, it is possible to affect users' behaviors indirectly via social influence.

However, recognizing influential users in a social network is difficult in practice. Although we can collect the topology of a social network and identify if a user is influential or not based on the number of her followers, it may contain false or weak edges that are ineffective at spreading influence [13]. A typical example is a user who has lots of fake followers on social media. Furthermore, toward distinct items or topics, the strength that an influential user influences her friends or followers could be diverse [14]. For instance, a famous blogger who focuses on commenting on movies may exert a stronger influence on her fans' choices of movies for watching but hardly affect music selections. In principle, we could conduct exhaustive surveys on every influential user to acquire the information of network topology and estimate the influence strength between each pair of users. Nevertheless, such surveys are very labor-intensive and impractical [15].

By considering the aforementioned challenging issues, is it possible to perform effective incentive allocation for realizing user incentivization in a social network, where the information about users and influence strength among users is absent, and the budget is restricted? To tackle this problem, we employ reinforcement learning (RL) to explore incentive allocation policies on a given network topology. More specifically, the RL agent is able to automatically adjust and determine the policy based on the interaction with users. Meanwhile, to better utilize the budget, the RL agent also learns the network representation in a way that is conducive to effective incentive allocation,

to explore influential users in the network and strategically incentivize them.

In our previous work [16], we conducted preliminary research work on tackling the incentive allocation problem in unknown social networks by using reinforcement learning. However, the details of how to apply reinforcement learning to the incentive allocation problem are not elaborated. In this paper, we first formulate the incentive allocation problem as a sequential decision problem. Then we also propose a novel end-to-end reinforcement learning-based framework, named Geometric Actor-Critic (GAC), which can learn to extract information from the network for generating incentive allocation policies. GAC can learn global representations for the entire network and local representations for each user. This allows the RL agent to estimate if the user is influential in the network, such that it determines the incentive value accordingly. In the experiments, three real-world datasets are deployed to construct the simulation environment for the evaluation. By comparing against existing incentive allocation approaches, the experimental results demonstrate the effectiveness of the proposed GAC.

The major contributions of this paper can be summarized as follow:

- We model the incentive allocation problem as a sequential decision problem by using Markov Decision Process (MDP), which makes using reinforcement learning to tackle the incentive allocation problem possible.
- We propose a novel reinforcement learning framework, i.e., Geometric Actor-Critic (GAC), to tackle the incentive allocation problem in unknown social networks. One key feature of GAC is that it only requires the basic information about the social network (i.e., the topology of the network) and the observation of users' behaviors for training and exploiting. Different from most existing incentive allocation approaches, the proposed RL approach is able to learn effective policies incentive allocation without knowing the attributes of users (e.g., users' preferences). To the best of our knowledge, this is the first work using reinforcement learning to tackle the incentive allocation problem in unknown social networks.
- We evaluate the proposed GAC on three real-world social network datasets. The experimental results demonstrate that GAC outperforms the compared approaches and prove its effectiveness.

The remainder of this paper is organized as follows: Section 2 reviews literature related to this research work. In Section 3, we first introduce the problem of incentive allocation in unknown networks. Then we present formal definitions and notations used in this paper. Meanwhile, we introduce the model used to build the simulation environment for Reinforcement learning. Subsequently, we give details about the proposed reinforcement learning-based framework in Section 4. Sections 5 and 6 demonstrate the experimental setup and results. The conclusion and future direction are illustrated in Section 7.

## 2. Related Work

In recent years, many studies have been devoted to designing effective incentive allocation approaches to incentivize users' behaviors. These studies aim to address the

problem of maximizing user incentivization under a budget constraint. Some incentive allocation approaches are designed for specific scenarios, where incentive values are determined by specific features or attributes, such as users’ preferences, location, and skill abilities [4, 17, 5, 18, 19, 20]. Although these approaches can perform effectively in incentive allocation, their performance could not be maintained when the information they rely on is unavailable. Meanwhile, the counterfactual inference-based approaches, which generate incentives by learning the logged feedback data from users, are also widely studied [21, 22]. A typical example is the budgeted multi-armed bandit problem, in which options are modeled as arms, and the target is to find out the most beneficial arm to select [23, 24, 25]. To apply such approaches in the real world, Singla et al. designed a UCB-based pricing mechanism, named DBP-UCB, to incentivize users of bike-sharing systems to help re-distribute bikes among stations in the real world [6]. However, these approaches assumed that users’ behaviors are independent but ignored the interaction and influence among users.

As aforementioned, a limited budget would restrict the number of users who can receive sufficient incentives and then be incentivized. To incentivize more users with a given budget, many studies have considered utilizing social influence which naturally exists among users to indirectly incentivize users [7, 1, 2]. These works assumed that the information about the social network including the network topology and the influence strength of each pair of users is known. Based on this assumption, these approaches can accurately locate influential users and incentivize them in priority. However, in the practical application, it is difficult to harvest complete information about the social network [13], and thus identifying influential users is challenging. Several studies have focused on a similar problem, called influence maximization in unknown networks, where the network structural information is incomplete or lacking, and the objective is to identify a set of users who can maximize influence diffusion [26, 27, 28]. The difference between theirs and our problem is that they do not consider the cost of incentivizing each candidate user while we need to do.

Many reinforcement learning models have emerged in the past years, such as DQN, A3C, DDPG, and PPO [29]. These approaches are designed for training an RL agent to automatically learn to achieve specific goals in an unfamiliar environment, e.g., learn to play a range of Atari 2600 video games at a superhuman level. Recently, some studies employ reinforcement learning to solve combinatorial problems on graphs, such as Travelling Salesman Problem [30], Online Vehicle Routing problem [31], and graph sampling problem [32]. Different from these works, our model aims to tackle the incentive allocation problem in an environment where only the topology of the social network is given. Our objective is to train an effective incentive allocation policy that can generate incentives for every user under a budget constraint. To the best of our knowledge, this is the first work using reinforcement learning to tackle the incentive allocation problem in unknown networks.

### 3. Preliminaries

#### 3.1. Problem Description

User incentivization is a complex problem, as the user’s decision-making can be affected by external factors, such as incentives. To effectively incentivize users’ be-

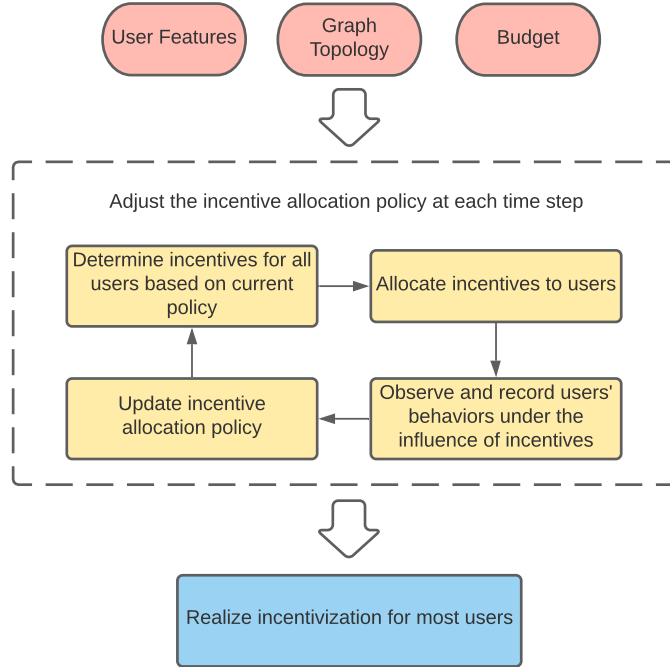


Figure 1: Learning policies for incentive allocation

haviors, applying a proper incentive allocation policy is essential. In this paper, the problem we focus on is to incentivize as many as possible users in a social network under a limited budget restriction. The information about the network is the topology only and no prior knowledge about users is given. To tackle this problem, it is necessary to learn the right policy for incentive allocation, which considers users' attitudes toward incentives and their influential roles in the social network simultaneously. The goal is that the approach is eventually able to incentivize users with appropriate incentives, and tends to incentivize influential users in priority.

Figure 1 shows the process of general incentive allocation. Given the topology of the network and a fixed amount of budget per time step, the system allocates incentives to all users based on the current allocation policy at every time step. At each time step, the approach would update its parameters based on the observation of users' behaviors, i.e., if the user accepts the incentive or not. This whole process repeats for a finite number of steps. After the last time step, the number of incentivized users reflects the performance of the approach.

### 3.2. Formal Definitions

In this section, we introduce definitions and notations used in this paper. Notations are listed in Table 1.

Table 1: Notations used in this paper

Notation	Description
$v_i$	A user $v_i$ in a social network
$b_{v_i,t}$	$v_i$ 's behavior at time step $t$
$G$	A directed graph of a social network
$V$	A set of users in $G$
$e_{ij}$	Influential relationship between $v_i$ and $v_j$
$E$	A set of edges in $G$
$w_{ij}$	Influence strength associated on $e_{ij}$
$N_{v_i}^{in}$	A set of $v_i$ 's one-hop in-neighbors
$N_{v_i}^{out}$	A set of $v_i$ 's one-hop out-neighbors
$z_m$	A behavior option
$z^*$	The behavior that the incentive provider expects
$Z$	A set of all behavior options
$p_{v_i,z_m}$	$v_i$ 's preference toward $z_m$
$o_{v_i,t}$	An incentive provided to $v_i$ at time step $t$
$B_t$	Remaining budget at time step $t$
$A^{in}$	The in-adjacency matrix
$A^{out}$	The out-adjacency matrix
$F$	The feature matrix
$f_{v_i}$	The feature vector for $v_i$
$u_t(v_i, z_m)$	$v_i$ 's utility toward $z_m$ at time step $t$
$k_t(v_i, z_m)$	Influence from $v_i$ 's in-neighbors toward $z_m$ at time step $t$
$S_t$	State at time step $t$
$\mathbf{a}_t$	The vector of an incentive action for all users at time step $t$
$R_t$	Step reward at time step $t$
$r_{v_i,t}$	Intermediate reward generated by $v_i$

**Definition 1.** An autonomous user agent representing a human user in a social network is defined as  $v_i$ .  $b_{v_i,t}$  denotes  $v_i$ 's behavior at time step  $t$ .

**Definition 2.** A social network is represented as a directed graph  $G = (V, E)$ , where  $V = \{v_1, \dots, v_i\}$  denotes a set of users, and  $E = \{e_{ij} | \{v_i, v_j\} \subseteq V, i \neq j\}$  denotes a set of edges in the network. In  $G$ ,  $e_{ij}$  represents an influential relationship between  $v_i$  and  $v_j$ . It implies that  $v_i$  can influence  $v_j$ 's behavior, and the influence strength associated on  $e_{ij}$  is represented as  $w_{ij} \in (0, 1]$ . We use  $N_{v_i}^{in}$  to represent a set of  $v_i$ 's one-hop in-neighbors who can influence  $v_i$  directly. Similarly,  $N_{v_i}^{out}$  represents a set of  $v_i$ 's one-hop out-neighbors who are directly influenced by  $v_i$ .

**Definition 3.** A behavior option  $z_m \in Z$  represents a specific behavior that  $v_i$  could choose, where  $Z = \{z_1, z_2, \dots, z_m\}$  denotes a finite option set. For each  $z_m$ ,  $v_i$  has a personal preference  $p_{v_i,z_m} \in [0, 1]$ , where higher  $p_{v_i,z_m}$  implies that  $v_i$  prefers to choose  $z_m$ , whereas smaller  $p_{v_i,z_m}$  implies conversely. To clarify, we define the behavior that the incentive provider expects  $v_i$  to select as  $z^*$ .

**Definition 4.** An incentive  $o_{v_i,t}$  denotes the incentive provided to user  $v_i$  from an incentive provider at time step  $t$ . The range of  $o_{v_i,t}$  can be  $[0, 1]$ . Note that the value of  $o_{v_i,t}$  is strictly constrained by the remaining budget  $B_t$ , i.e.,  $o_{v_i,t} \leq B_t$ .

Furthermore, there are several matrices that would be used in Section 4. By default, we use bold uppercase and lowercase letters, e.g.,  $\mathbf{A}$  and  $\mathbf{x}$ , to represent matrices and vectors, respectively. Let  $\mathbf{A}^{out} \in \mathbb{R}^{|V| \times |V|}$  be out-adjacency matrix, where each entry  $\mathbf{A}^{out}[i, j] = 1$  if  $e_{ij} \in E$ . Similarly, we use  $\mathbf{A}^{in} \in \mathbb{R}^{|V| \times |V|}$  to represent in-adjacency matrix, where each entry  $\mathbf{A}^{in}[j, i] = 1$  if  $e_{ji} \in E$ . The features matrix of all users is represented by  $\mathbf{F}$ , where each row represents a specific user's feature  $\mathbf{f}_{v_i} = [o_{v_i,t}, b_{v_i,t}]$ . In this paper, we use the one-hot representation for  $b_{v_i,t}$ .

### 3.3. Simulation of Environment

To model the process of user incentivization, we build a simulation-based environment by adopting the Agent-based Decision-making Model (ADM) [33], in which two types of agents are involved, i.e., the RL agent and user agents. The RL agent stands for the incentive provider and generates incentives for users. While user agents represent human users in the real world. All user agents are required to behave at every time step, and each user agent always chooses an item with the highest user utility. Equation 1 formulates this behavior rule, where  $u_t(v_i, z_m)$  denotes the user utility of  $z_m$ .

$$b_{v_i,t} = \arg \max_{z_m \in Z} u_t(v_i, z_m) \quad (1)$$

The user utility  $u_t(v_i, z_m)$  is formulated in Equation 2, which consists of the user's preferences  $p_{v_i, z_m}$ , received incentive  $o_{v_i,t}$ , and the social influence  $k_t(v_i, z_m)$  from its neighbors simultaneously. Namely, users select items purely based on their own preferences and social influence if no incentive is provided.

$$u_t(v_i, z_m) = \begin{cases} p_{v_i, z_m} + k_t(v_i, z_m) + o_{v_i,t}, & z_m = z^* \\ p_{v_i, z_m} + k_t(v_i, z_m), & otherwise \end{cases} \quad (2)$$

In the ADM,  $b_{v_i,t}$  can produce influence on  $b_{v_j,t+1}$ ,  $\forall v_j \in N_{v_i}^{out}$ . Hence, multiple influences for supporting different items may co-exist in the meanwhile, which could form the conflict in affecting the focal user's behavior. Equation 3 describes how the ADM formulates the social influence exerted on  $v_i$  at time step  $t$ , where  $w_{ji}$  denotes the strength of influence associated with edge  $e_{ji}$ .

$$k_t(v_i, z_m) = \sum_{\substack{v_j \in N_{v_i}^{in}, \\ b_{v_j,t-1} = z_m}} w_{ji} \quad (3)$$

Although the ADM requires knowledge about users' preferences and the strength of influence associated with all edges to conduct the simulation, the proposed algorithm does not require such prior knowledge.

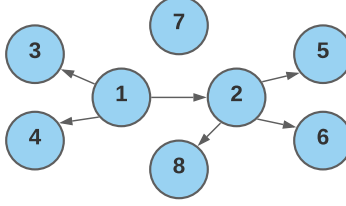


Figure 2: An example of a simple social network

## 4. Methodology

In this section, we introduce the proposed approach for training an agent that can automatically produce effective incentive allocation policies. We first formulate the incentive allocation problem as a Markov Decision Process (MDP) problem, with a carefully designed reward function to enable the training process can fit different networks and budget restrictions. Subsequently, we introduce the structure of the proposed GAC, and how we encode local and global embeddings as the state, which provides contextual information to the agent to determine the incentive allocation policy. Last, we introduce the process for training the proposed approach.

### 4.1. Markov Decision Process Formulation for Incentive Allocation Problem

We first formulate the incentive allocation problem as a MDP. As aforementioned, the process of incentive allocation goes through a range of time steps before the finish. At each time step, the RL agent generates incentives for all users based on the current state of the network. After assigning incentives to all users, users would decide to accept the incentive or not and return the feedback to the RL agent. The state of the network subsequently changes at the end of the time step. Hence, the incentive allocation problem can be modeled as a sequential decision problem. In this MDP, the corresponding formal representation of transition at each time step can be formulated as  $(s_t, \mathbf{a}_t, s_{t+1}, R_t)$ , where  $s_t$  denotes the current state of the network,  $\mathbf{a}_t$  represents the action, i.e., incentives allocated to all users,  $s_{t+1}$  is the next state after providing incentives, and  $R_t$  is the step reward.

Notably, action  $\mathbf{a}_t$  is a vector where each element is a continuous value representing an incentive for all users. Hence, before assigning incentives to users, we need to map  $\mathbf{a}_t$  to an independent incentive for each user. The step reward  $R_t$  is a cumulative reward which is formulated in Equation 4, where  $r_{v_i,t}$  denotes the intermediate reward for assigning an incentive to  $v_i$ . Step reward is only available at the end of each time step.

$$R_t = \sum_{v_i \in V} r_{v_i,t} \quad (4)$$

An intermediate reward  $r_{v_i,t}$  is calculated by using Equation 5. In the equation,  $\alpha$  is the indicator of the user activation.  $\alpha = 1$  if user  $v_i$  accepts the incentive and choose  $z^*$ ,



and  $\alpha = -1$  if conversely. This equation consists of two parts. The former part indicates the importance of  $v_i$  in a social network, where  $|N_{v_i}^{out}|$  and  $|N_{v_i}^{in}|$  represent the number of neighbors who can be influenced by  $v_i$  and who can influence  $v_i$ , respectively. To better explain the former part, we take the social network in Figure 2 as an example. In this network, each node represents a specific user, and each directed edge represents the influence direction. Empirically, if a user can influence more users in the network, we can regard this user is more valuable to be incentivized. However, in spite that users 1 and 2 can both influence three users, incentivizing user 1 first could be a better choice when the budget is insufficient for incentivizing both of them, as user 1 can influence user 2. Based on this rule, the importance of each user can be calculated and distinguished. The latter part represents the efficiency of spending the budget. If a user can be incentivized with fewer incentives, then the budget can be saved for incentivizing subsequent users. Different from step reward, the intermediate reward would be generated once  $v_i$  responds to the provided incentive.

$$r_{v_i,t} = \alpha \left( 1 + \frac{|N_{v_i}^{out}| - |N_{v_i}^{in}|}{|V|} \right) + \frac{\alpha + 1}{2} \cdot \frac{B - o_{v_i,t}}{B} \quad (5)$$

Algorithm 1 describes how the RL agent incentivizes users at every time step. The input is an action  $\mathbf{a}_t$  generated by the RL agent, and a fixed budget  $B$ , and the output is step reward  $R_t$  and the log  $L_t$ , which is used to store records of users’ behaviors. Line 1 sets the remaining budget as  $B$  and step reward as 0, and Line 2 initializes  $L_t$  as an empty list. Lines 4-6 first map the value of incentive for  $v_i$  from  $\mathbf{a}_t$ , and then check if the remaining budget  $B_t$  is sufficient for allocating the incentive. In Lines 7-8,  $v_i$  makes a decision and takes the behavior based on its user utilities. Subsequently, variables are updated based on the user’s behavior  $b_{v_i,t}$  in Lines 9-13. At last, in Lines 14-15, intermediate reward and step reward are calculated, and the record of  $v_i$ ’s behavior is added to  $L_t$  in Line 16.

#### 4.2. Geometric Actor-Critic

To tackle the incentive allocation problem in a social network, it is essential to take both the structure of the network and users’ features into consideration when generating incentive policies. Due to the variable size and complexity of a social network, it is also necessary to encode the information of a network in a low-level graph representation.

Figure 3 shows the structure of the proposed GAC. The input of GAC includes the user feature matrix and two different adjacency matrices, i.e., the in-adjacency matrix and the out-adjacency matrix. The reason we input two adjacency matrices is that a social network is typically a directed graph structure, and its adjacency matrix is not symmetric. To keep all structural information about the graph, we deploy two independent components consisting of Graph Neural Networks (GNNs) to encode both the out-adjacency  $\mathbf{A}^{out}$  and the in-adjacency matrices  $\mathbf{A}^{in}$ , respectively. In each network encoding component, a GraphSage [34] is deployed first to learn refined user features by aggregating features from their neighbors. In this work, we use the Mean “variant” of GraphSage by default. By taking the refining user features from the user’s in-neighbors as an example, the formulation of a GraphSage layer can be formulated in Equations 6 and 7. In this aggregation,  $\mathbf{h}_{N_{v_i}^{in}}^q$  representing the vector representation of  $N_{v_i}^{in}$  is first

---

**Algorithm 1:** Assign Incentives to users at time step  $t$ 

---

**Input:** Action  $\mathbf{a}_t$ , budget  $B$   
**Output:** Step reward  $R_t$ , log  $L_t$

- 1  $B_t := B; R_t := 0;$
- 2 Initialize  $L_t$  as an empty list;
- 3 **for**  $v_i \in V$  **do**
- 4     Map  $o_{v_i,t}$  from  $\mathbf{a}_t$ ;
- 5     **if**  $o_{v_i,t} > B_t$  **then**
- 6          $o_{v_i,t} = B_t;$
- 7     Calculate  $u_t(v_i, z_m), \forall z_m \in Z$  by using Equation 2;
- 8      $b_{v_i,t} = \arg \max_{z_m \in Z} u_t(v_i, z_m);$
- 9     **if**  $b_{v_i,t} == z^*$  **then**
- 10          $B_t = B_t - o_{v_i,t};$
- 11          $\alpha := 1;$
- 12     **else**
- 13          $\alpha := -1;$
- 14     Calculate  $r_{v_i,t}$  by using Equation 5;
- 15      $R_t = R_t + r_{v_i,t};$
- 16     Add  $(v_i, o_{v_i,t}, b_{v_i,t})$  to  $L_t$ ;
- 17 **Return**  $R_t, L_t$

---

generated by aggregating representation of all  $v_i$ 's in-neighbors. Then, the layer concatenates  $v_i$ 's representation  $\mathbf{h}_{v_i}^{q-1}$  and  $\mathbf{h}_{N_{v_i}^{in}}^q$ . Subsequently, the concatenated vector is fed through a fully connected layer with a weight matrix  $\mathbf{W}^q$  and a nonlinear activation function  $ReLU$  to transform the representation of the user  $v_i$ . In this equation,  $q$  denotes  $q$ -hop neighbors considered in the GraphSage. For example, when  $q = 1$ , it means we only consider aggregating features from  $v_i$ 's one-hop neighbors, and then  $\mathbf{h}_{v_i}^{q-1}$  becomes  $\mathbf{h}_{v_i}^0$  representing the input user's feature. After the first GraphSage layer, we can obtain local representation for all users, i.e., node embeddings.

$$\mathbf{h}_{N_{v_i}^{in}}^q \leftarrow MEAN_q(\{\mathbf{h}_{v_j}^{q-1}, v_j \in N_{v_i}^{in}\}) \quad (6)$$

$$\mathbf{h}_{v_i}^q \leftarrow ReLU(\mathbf{W}^q \cdot CONCAT(\mathbf{h}_{v_i}^{q-1}, \mathbf{h}_{N_{v_i}^{in}}^q)) \quad (7)$$

After learning the node embeddings for all users, we use DIFFPOOL [35] to learn the global representation of the entire social network, which aggregates user features in a hierarchical manner. DIFFPOOL can learn hierarchical representations of the input network by iteratively coarsening the network, mapping users to a set of clusters. Specifically, DIFFPOOL first generates the embedding matrix  $\mathbf{Y}^l \in \mathbb{R}^{n_l \times d}$  and the assignment matrix  $\mathbf{V}^l \in \mathbb{R}^{n_l \times n_{l+1}}$ , where  $l$  denotes the index of DIFFPOOL,  $n_l$  denotes the number of clusters or users in the input coarsened network,  $n_{l+1}$  denotes the number of clusters in the output coarsened network, and  $d$  denotes the dimension of embeddings. Equations 8 and 9 formulate how to generate  $\mathbf{Y}^l$  and  $\mathbf{V}^l$  via two independent GNNs,

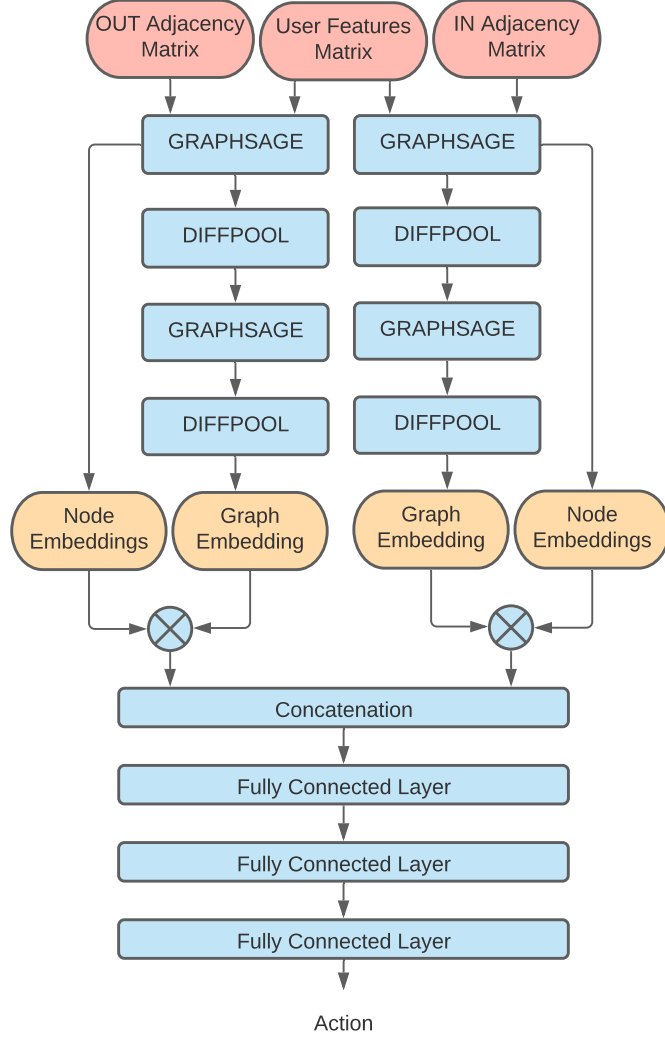


Figure 3: Architecture of GAC

where  $A^l$  and  $X^l$  represent the corresponding adjacency matrix and node embeddings matrix. Here, we use two GraphSages for GNNs. Note that when  $l = 0$ ,  $A^l$  is the original adjacency matrix, and  $X^l$  denotes the node embeddings generated by the first GraphSage.

$$Y^l = GNN_{embed}^l(A^l, X^l) \quad (8)$$

$$V^l = GNN_{pool}^l(A^l, X^l) \quad (9)$$

Then, we can obtain a new embeddings matrix  $X^{l+1} \in \mathbb{R}^{n_{l+1} \times d}$  for each cluster in the coarsened network and a new coarsened adjacency matrix  $A^{l+1} \in \mathbb{R}^{n_{l+1} \times n_{l+1}}$  by using Equations 10 and 11. The output  $X^{l+1}$  and  $A^{l+1}$  would be fed to the next GraphSage and then DIFFPOOL.

$$X^{l+1} = V^{lT} Y^l \quad (10)$$

$$A^{l+1} = V^{lT} Y^l V^l \quad (11)$$

At last, the coarsened network would only have one cluster, which is a vector representation for the entire network. In this study, we use the first DIFFPOOL to map all users into 16 clusters, and the second DIFFPOOL to generate the global representation for the corresponding network.

Once the local and global representations are obtained, i.e., node embeddings for all users and graph embedding, they can be combined via a matrix-vector product as described in Equations 12 and 13. In these two equations,  $\mathbf{x}_{in} \in \mathbb{R}^{1 \times d}$  and  $\mathbf{x}_{out} \in \mathbb{R}^{1 \times d}$  represent graph embeddings generated by the in-adjacency and the out-adjacency matrices, respectively. While  $\mathbf{H}_{in} \in \mathbb{R}^{|V| \times d}$  and  $\mathbf{H}_{out} \in \mathbb{R}^{|V| \times d}$  represent the corresponding node embeddings.  $d$  denotes the dimensional length of embedding. The matrices of node embeddings need to be transposed to make the calculation feasible. Through Equations 12 and 13, two new vectors  $\mathbf{x}_{in}^* \in \mathbb{R}^{1 \times |V|}$  and  $\mathbf{x}_{out}^* \in \mathbb{R}^{1 \times |V|}$  which encode local and global representations together can be generated.

$$\mathbf{x}_{in}^* = \mathbf{x}_{in} \mathbf{H}_{in}^T \quad (12)$$

$$\mathbf{x}_{out}^* = \mathbf{x}_{out} \mathbf{H}_{out}^T \quad (13)$$

The GAC then concatenates  $\mathbf{x}_{in}^*$  and  $\mathbf{x}_{out}^*$  together. The concatenated vector would be normalized by using L2-Norm and subsequently fed through three fully connected layers with nonlinear activation functions. The first two layers use *ReLU* as activation functions, while the last layer use *tanh* as the activation function. The output  $\mathbf{a}_t \in \mathbb{R}^{1 \times |V|}$  is the action representation, where the range of each entry is from -1 to 1. In this case, we need to re-scale  $\mathbf{a}_t$  to the range from 0 to 1 before allocating incentives for users.

### 4.3. Model Training

Algorithm 2 shows the training process for the proposed GAC. Before starting the training process, we initialize actor and critic networks of GAC as well as the replay buffer  $\mathcal{B}$ . At the beginning of every episode, we assign no incentive to all users and obtain the observation of their behaviors to generate node features matrix  $\mathbf{F}$ .  $\omega$  in Lines 6 and 14 represents the ratio of engaged users in the network. The state is represented by  $S = (A^{out}, A^{in}, \mathbf{F})$ . In the initial *EXP* episodes, the action is captured from the normal distribution  $\mathcal{N}(0, 1)$  for pure exploration. After *EXP* episodes,  $\mathbf{a}_t$  is generated by GAC directly. To keep the approach exploring, we add noise  $\epsilon$  on the generated  $\mathbf{a}_t$ , where  $\epsilon \sim \mathcal{N}(-\omega, 1)$ . The reason we capture noise from this distribution is that we want to decrease the incentive allocated to users when most users in the network have been incentivized. At that moment, most users would be exerted influence from

---

**Algorithm 2: GAC Training**

---

**Input:** Number of episodes  $EP$ , number of time steps  $T$ , number of episodes for exploration  $EXP$ , budget  $B$ , in-adjacency matrix  $A^{in}$ , out-adjacency matrix  $A^{out}$

```
1 Initialize actor network  $\pi_\phi$ , and critic networks  $Q_{\theta_1}$  and  $Q_{\theta_2}$  with random
  parameters  $\phi, \theta_1, \theta_2$ ;
2 Initialize target networks  $\phi' \leftarrow \phi, \theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2$ ;
3 Initialize replay buffer  $\mathcal{B}$ ;
4 for  $episode = 1$  to  $EP$  do
5   Incentivize all users with no incentive and observe users' behavior to
   initialize feature matrix  $F$ ;
6    $\omega := \frac{|\{v_i | v_i \in V, b_{v_i,0} = z^*\}|}{|V|}$ ;
7   for  $t = 1$  to  $T$  do
8      $S = (A^{out}, A^{in}, F)$ ;
9     if  $episode \leq EXP$  then
10       $\mathbf{a}_t \sim \mathcal{N}(0, 1)$ ;
11    else
12       $\mathbf{a}_t \sim \pi_\phi(S) + \epsilon, \epsilon \sim \mathcal{N}(-\omega, 1)$ ;
13       $R_t, L_t = AssignIncentive(\mathbf{a}_t, B)$  (Algorithm 1);
14       $\omega := \frac{|\{v_i | v_i \in V, b_{v_i,t} = z^*\}|}{|V|}$ ;
15      Observe step reward  $R_t$  and create new Feature matrix  $F'$  based on  $L_t$ ;
16      Store transition tuple  $(F, \mathbf{a}_t, R_t, F')$  in  $\mathcal{B}$ ;
17       $F = F'$ ;
18    if  $episode > EXP$  then
19      Sample mini-batch of transitions from  $\mathcal{B}$  and update actor and
      critic networks;
20      Update target networks by using soft-replacement;
```

---

neighbors affecting them to choose  $z^*$ , and then they might be indirectly incentivized by their neighbors. In Line 13, the generated action  $\mathbf{a}_t$  and budget  $B$  would be inputted to the environment which we introduce in Algorithm 1. The environment returns the step reward  $R_t$  and users' behavior  $\log L_t$ . The received  $L_t$  is used to update  $\omega$  and create a new user features matrix. At every time step, we store the experience tuple  $(F, \mathbf{a}_t, R_t, F')$  in the replay buffer  $\mathcal{B}$ . To save memory, we do not save the adjacency matrices to the replay buffer, as the network topology would not change in this study.

After  $EXP$  episodes, we can sample experiences from the replay buffer to train parameters in both actor and critic networks of GAC. Inspired by TD3 [36], which is a variant of the popular DDPG algorithm [37] for continuous control, we also create two critic networks in the GAC (Line 2). The major modification of TD3 is using a Double Q network and adding noise to the policy when computing the target value. At last, the target networks are updated by using soft-replacement.

Table 2: Statistics of datasets

Dataset	$ V $	$ E $	Avg. Degree
Dolphins	62	159	5.1
Twitter	236	2,478	21.0
Wiki-Vote	889	2,914	6.6

## 5. Experimental Setup

### 5.1. Data Preparation

To evaluate the performance of the proposed GAC, the following three datasets are used to deploy social networks:

- **Dolphins**<sup>1</sup> dataset represents a social network of bottlenose dolphins, where a node represents a dolphin, and an edge represents frequent associations between dolphins [38]. This network contains 62 nodes and 159 edges in total, and the average degree for each node is 5.1.
- **Twitter**<sup>2</sup> dataset contains 973 directed networks, 81,306 users and 1,768,149 edges in total [39]. To diminish the running time, a sub-network that contains 236 users and 2,478 edges is selected. The average degree of the sub-network is 21.0.
- **Wiki-Vote**<sup>3</sup> dataset contains all the voting data from the inception of Wikipedia till January 2008 [40]. This dataset contains 889 nodes and 2,914 edges, where each directed edge from node  $i$  to node  $j$  represents that user  $i$  voted user  $j$ . The average degree is 6.6.

To simulate user agents' behaviors in the ADM, we also assign random preferences  $p_{v_i, z_m}, \forall z_m \in Z$  for all users in these three datasets. We suppose that four behavior options can be selected, and  $z_0$  is regarded as the expected action  $z^*$ . For the strength of influence associated with each edge, we also assign a random value from 0 to 1, and ensure that the sum of influence strength from each user's in-neighbors would not exceed 1, i.e.,  $\sum_{v_j \in N_{v_i}^{in}} w_{ji} \leq 1, \forall v_i \in V$ . The statistics of three datasets are listed in Table 2.

### 5.2. Baseline Approaches

The performance of the proposed GAC is evaluated by comparing it with the following approaches:

- **No Incentive** approach implies that no incentive would be allocated to users. Hence, all users would make decisions only based on their preferences and social influence from neighbors.

<sup>1</sup><http://networkrepository.com/soc-dolphins.php>

<sup>2</sup><https://snap.stanford.edu/data/ego-Twitter.html>

<sup>3</sup><http://networkrepository.com/soc-wiki-Vote.php>

- **Uniform Allocation** is a Naïve approach, which allocates the same incentives to users. The incentive values are determined by the number of users and the budget amount.
- **DGIA-IPE** is an adaptive incentive allocation approach which consists of two components, i.e., DGIA and IPE [33]. Based on the observation of users’ behaviors, DGIA can estimate the users’ sensitivity toward incentives, and IPE can estimate influential relationships among users. Finally, DGIA-IPE determines the incentive values based on the results of estimation for each user.
- **DBP-UCB** is a dynamic pricing algorithm proposed by Singla et al. for engaging users in the system to participate in the bike re-positioning process [6]. Different from DGIA and the proposed GAC, DBP-UCB models each price option as a discrete option. The parameters of DBP-UCB would be updated based on the observation of users’ behaviors.

In the experiments, the default metric used to evaluate the performance of approaches is the number of users who take  $z^*$ .

### 5.3. System Settings

We conduct simulation-based experiments to evaluate the proposed approach, i.e., GAC, by comparing it with the other four baseline approaches. For GAC and its variants, we set the number of units in GraphSage layers and fully connected layers as 32 and 64, respectively. We use Adam Gradient Descent [41] to optimize the model, and the learning rates for Actor and Critic networks are  $3e-4$  and  $3e-5$ , respectively. The size of the replay buffer is set as  $5e4$ , and the batch size of samples is 200. The GAC and its variants are trained by  $1e4$  episodes, with 10 time steps in every episode by default. The initial 1500 episodes are used for pure exploration, where the action is generated by using normal distribution  $\mathcal{N}(0, 1)$ . At last, the model parameters which can generate the most optimal incentive allocation policy would be kept. We set up DGIA-IPE by using the same parameters used in the original paper. For DBP-UCB, we give ten price options from 0 to 1 with an interval of 0.1. In the evaluation, the simulation for each approach lasts 150 time steps, and the budget will be refilled at the beginning of each time step.

## 6. Experimental Results

### 6.1. Impact of GAC’s Structure

In Section 4, we explain that the purpose of inputting both in-adjacency and out-adjacency matrices is to prevent losing important information on the network. To prove the importance of this operation, we compare the GAC with its two simplified versions, i.e., GAC-IN and GAC-OUT. The difference between these two simplified GACs is the input matrices. Besides the feature matrix, GAC-IN only takes the in-adjacency matrix as input, while GAC-OUT requires the out-adjacency matrix.

We compare the performance of these three variants by using the Dolphins and Twitter datasets, respectively. As we can observe from Figure 4, GAC slightly outperforms the simplified GAC in both datasets. Although the performance of GAC-OUT

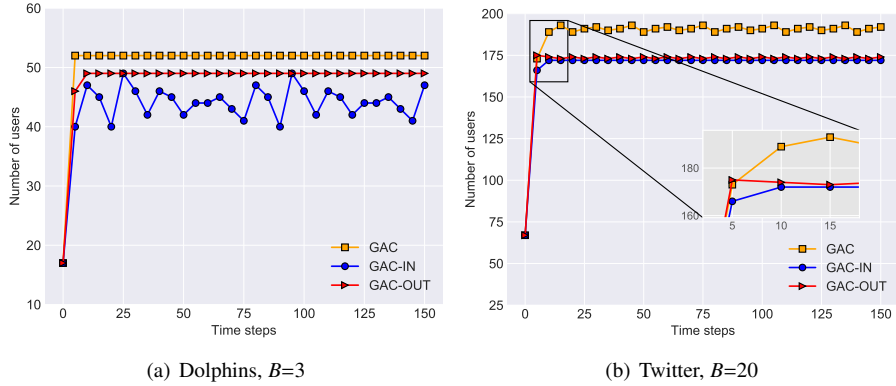


Figure 4: Performance comparison of GAC and its variants

is slightly worse than GAC, it eventually converges to a stable stage. By contrast, GAC-IN performs the worst in the Dolphins dataset. A possible reason is that GAC-IN only considers users' in-neighbors when generating incentives, such that the influential users in the network are difficult to be identified. However, GAC-IN performs similarly as GAC-OUT in the Twitter dataset. In spite that GAC-OUT reaches the convergence faster than GAC-IN, the gap between them can be ignored after convergence. It is due to that the average degree in the Twitter network is higher than the Dolphins network, and the users in the Twitter network who can influence others might also be influenced by their neighbors.

## 6.2. Impact of Noise Distribution

Many algorithms for continuous control, such as TD3, add noise captured from a normal distribution  $\mathcal{N}(\mu, \sigma^2)$  to the policy, where  $\mu$  is typically set as 0, and  $\sigma$  is a fixed value. By applying this noise, the algorithm is able to explore different policies and eventually obtains success. However, we found that setting  $\mu$  as 0 cannot help to explore effective policy in the incentive allocation problem where social influence exists. The reason is that, with the increasing number of users who choose  $z^*$  in the network, more users start affecting their neighbors to select  $z^*$ . At that moment, it could be unnecessary to allocate incentives to some users, as they can be incentivized by the influence generated by their neighbors.

Hence, in this study, we consider capturing noise from  $\mathcal{N}(-\omega, 1)$ , where  $\omega$  represents the ratio of users who select  $z^*$  in the network. We compare the performance of GAC by using three different noise distributions, i.e.,  $\mathcal{N}(-\omega, 1)$ ,  $\mathcal{N}(0, 0.2)$ , and  $\mathcal{N}(0, 1)$ . As we can observe from Figures 5(a) and 5(b), using  $\mathcal{N}(-\omega, 1)$  can attempt diverse incentive allocation policies and receive different cumulative step rewards from the environment. These experiences are beneficial for GAC to learn if the policy is good or not. By contrast, using  $\mathcal{N}(0, 0.2)$  or  $\mathcal{N}(0, 1)$  to generate noise makes the exploration very limited, leading the learned experience might be very similar.

The results from Figures 5(c) and 5(b) also demonstrate that capturing noise from  $\mathcal{N}(-\omega, 1)$  is more beneficial for training GAC, as it allows GAC to explore diverse



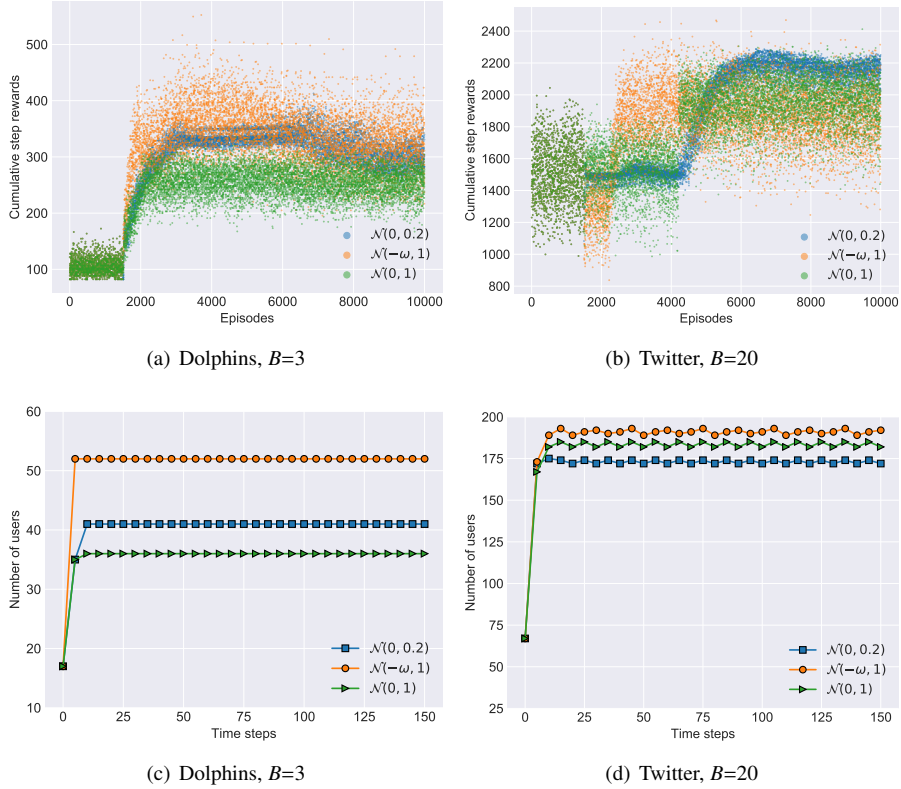
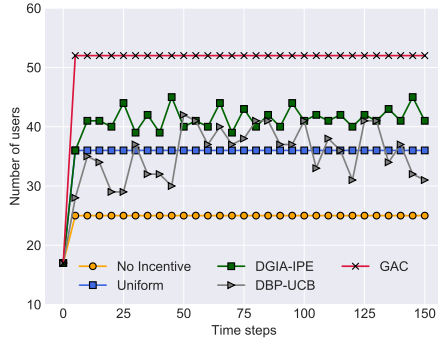


Figure 5: Performance comparison of adding different noise

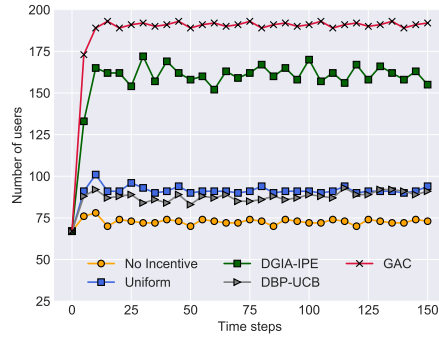
policies, which makes the model be able to learn from both good and bad experiences. By contrast, using  $\mathcal{N}(0, 0.2)$  or  $\mathcal{N}(0, 1)$  to generate noise makes GAC easily learn the sub-optimal policies, but fail to learn the optimal policy. Notably,  $\mathcal{N}(0, 1)$  outperforms  $\mathcal{N}(0, 0.2)$  in the Twitter network, while it underperforms  $\mathcal{N}(0, 0.2)$  in the Dolphins network. A possible reason is that social influence plays a more crucial role in the Twitter network compared to the Dolphins network. Meanwhile, making  $\sigma^2$  larger in the normal distribution implies that the probability of capturing a larger noise increases, such that it is possible to explore the policies that allocate users no incentive.

### 6.3. Comparison with Baseline Approaches

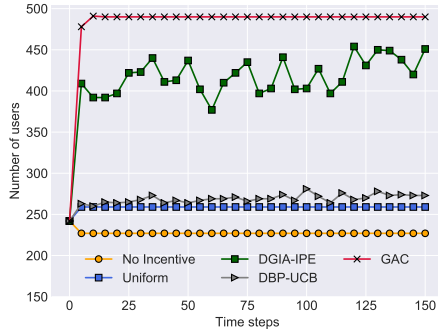
In this experiment, we compare the performance of the proposed GAC with the other four baseline approaches in three different social network datasets. As we can observe from Figure 6(a), the proposed GAC is able to incentivize more users than the other baseline approaches in the Dolphins network. The performance of DGIA-IPE is slightly better than Uniform allocation and DBP-UCB. The reason is that the IPE algorithm would estimate the influential relationships among users while Uniform allocation and DBP-UCB cannot.



(a) Dolphins,  $B=3$



(b) Twitter,  $B=20$



(c) Wiki-Vote,  $B=40$

Figure 6: Performance comparison of GAC and baseline approaches

In the Twitter network, GAC and DGIA-IPE produce much better performance than the other three baseline approaches, as they both consider social influence when generating incentive allocation policy. In comparison with DGIA-IPE, the proposed GAC converges rapidly and performs better. DBP-UCB has a very similar performance as the Uniform allocation. It implies that DBP-UCB is ineffective in incentivizing users in such a dense social network environment.

Figure 6(c) shows the performance of approaches in the Wiki-Vote network. The proposed GAC could still outperform other compared approaches. Similarly, DGIA-IPE yields a better performance than DBP-UCB and Uniform allocation, but underperforms GAC. Also, we notice that the performance of DGIA-IPE appears not very stable. A possible reason is that the average degree of Wiki-Vote is much lower than that of the Twitter network, and the IPE algorithm fails to estimate influence strength well in such a sparse social network. Meanwhile, different from its performance in the Twitter network, DBP-UCB slightly outperforms the Uniform allocation this time. This might also be caused by the sparsity of the network, as DBP-UCB does not take social influence into consideration when generating incentives.

#### 6.4. Discussion

In the experiment, we simulate a social environment and the process of incentivization in a social network, to evaluate the performance of the proposed GAC in solving the incentive allocation problem. The proposed GAC is evaluated by using three different real-world datasets. We first evaluate the impact of different architectures and noise distributions on the proposed GAC, and then compare the performance of GAC against existing approaches for the incentive allocation problem. Based on the experimental results, insights can be proposed as follows:

- The results from experiment 1 demonstrate that inputting the complete information of the network, i.e., in-adjacency and out-adjacency matrices, can produce better performance. Meanwhile, the out-adjacency matrix could be more important than the in-adjacency matrix in the GAC.
- Experiment 2 reveals that different noise distributions can cause impacts on the performance of incentive allocation. Specifically, GAC with a dynamic noise distribution outperforms that with a static noise distribution.
- The results from experiment 3 demonstrate that GAC outperforms the existing baseline approaches for the incentive allocation in unknown social networks. It reveals that learning representations for both users and the network can lead to a better result of incentive allocation.

### 7. Conclusion and Future Work

In this paper, we propose a Reinforcement Learning-based framework, named Geometric Actor-Critic (GAC), to solve the incentive allocation problem in unknown social networks with a budget limitation, where only the information about the network topological structure is available, and the influence strength and users' attributes are not provided. To tackle this problem, the proposed GAC learns to represent the network from both global and local perspectives, and generates incentive allocation policies based on learned information. The trained GAC is evaluated by comparing it with other baseline approaches in three different real-world social network datasets. The experimental results demonstrate that GAC outperforms other approaches in all three datasets under a budget limitation.

Although GAC lights a potential direction for the incentive allocation problem, the shortcomings of GAC cannot be ignored. As GAC requires two adjacency matrices and node features matrix as input, it spends a lot of memory space to temporarily store them. It also makes GAC impossible to handle extremely large social networks. In future work, we will keep investigations on effective Reinforcement Learning-based framework for the incentive allocation problem, and improve the proposed GAC.

### References

- [1] D. Zhao, B. Li, J. Xu, D. Hao, N. R. Jennings, Selling Multiple Items via Social Networks, in: Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems, 2018, pp. 68–76.

- [2] W. Zhang, D. Zhao, Y. Zhang, Incentivize Diffusion with Fair Rewards, in: Proceedings of the 24th European Conference on Artificial Intelligence, 2020, pp. 251–258.
- [3] Y. Singer, M. Mittal, Pricing Mechanisms for Crowdsourcing Markets, in: Proceedings of the 22nd international conference on World Wide Web, 2013, pp. 1157–1166.
- [4] X. Gan, X. Wang, W. Niu, G. Hang, X. Tian, X. Wang, J. Xu, Incentivize Multi-Class Crowd Labeling Under Budget Constraint, *IEEE Journal on Selected Areas in Communications* 35 (2017) 893–905.
- [5] C. Qiu, A. Squicciarini, B. Hanrahan, Incentivizing distributive fairness for crowdsourcing workers, in: Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, 2019, p. 404–412.
- [6] A. Singla, M. Santoni, G. Bartók, P. Mukerji, M. Meenen, A. Krause, Incentivizing Users for Balancing Bike Sharing Systems, in: Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015, pp. 723–729.
- [7] S. Wu, Q. Bai, S. Sengvong, GreenCommute: An Influence-Aware Persuasive Recommendation Approach for Public-Friendly Commute Options, *Journal of Systems Science and Systems Engineering* 27 (2018) 250–264.
- [8] A. Singla, A. Krause, Truthful Incentives in Crowdsourcing Tasks Using Regret Minimization Mechanisms, in: Proceedings of the 22nd international conference on World Wide Web, 2013, pp. 1167–1178.
- [9] Y. Wu, F. Li, L. Ma, Y. Xie, T. Li, Y. Wang, A Context-Aware Multiarmed Bandit Incentive Mechanism for Mobile Crowd Sensing Systems, *IEEE Internet of Things Journal* 6 (2019) 7648–7658.
- [10] Y. Li, J. Fan, Y. Wang, K.-L. Tan, Influence Maximization on Social Graphs: A Survey, *IEEE Transactions on Knowledge and Data Engineering* 30 (2018) 1852–1872.
- [11] M. J. Lovett, R. Peres, R. Shachar, On brands and word of mouth, *Journal of Marketing Research* 50 (2013) 427–444.
- [12] J. Axsen, C. Orlebar, S. Skippon, Social influence and consumer preference formation for pro-environmental technology: The case of a U.K. workplace electric-vehicle study, *Ecological Economics* 95 (2013) 96–107.
- [13] B. Wilder, N. Immerlica, E. Rice, M. Tambe, Maximizing Influence in an Unknown Social Network, in: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, 2018, pp. 4743–4750.
- [14] J. Tang, J. Sun, C. Wang, Z. Yang, Social Influence Analysis in Large-Scale Networks, in: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2009, pp. 807–816.

- [15] T. W. Valente, P. Pumpuang, Identifying opinion leaders to promote behavior change, *Health Education & Behavior* 34 (2007) 881–896.
- [16] S. Wu, Q. Bai, W. Li, Learning Policies for Effective Incentive Allocation in Unknown Social Networks, in: *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, 2021, pp. 1701–1703.
- [17] S. Sengvong, Q. Bai, Persuasive public-friendly route recommendation with flexible rewards, in: *2017 IEEE International Conference on Agents (ICA)*, 2017, pp. 109–114.
- [18] S. Xiao, L. Lv, L. Guo, Y. Chen, S. Yang, Z. Jiang, J. Zhu, Model-based constrained MDP for budget allocation in sequential incentive marketing, in: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 971–980.
- [19] J. Li, Y. Zhu, J. Yu, Redundancy-Aware and Budget-Feasible Incentive Mechanism in Crowd Sensing, *The Computer Journal* 63 (2019) 66–79.
- [20] S. Wu, W. Li, H. Shen, Q. Bai, Identifying influential users in unknown social networks for adaptive incentive allocation under budget restriction, 2021. [arXiv:2107.05992](https://arxiv.org/abs/2107.05992).
- [21] S. Wu, Q. Bai, Incentivizing Long-Term Engagement Under Limited Budget, in: *PRICAI 2019: Trends in Artificial Intelligence*, 2019, pp. 662–674.
- [22] R. Lopez, C. Li, X. Yan, J. Xiong, M. I. Jordan, Y. Qi, L. Song, Cost-Effective Incentive Allocation via Structured Counterfactual Inference, in: *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020, pp. 4997–5004.
- [23] L. Tran-Thanh, A. Chapman, E. M. de Cote, A. Rogers, N. R. Jennings, Epsilon-First Policies for Budget-Limited Multi-Armed Bandits, in: *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2010, pp. 1211–1216.
- [24] L. Tran-Thanh, A. Chapman, A. Rogers, N. R. Jennings, Knapsack Based Optimal Policies for Budget-Limited Multi-Armed Bandits, in: *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012, pp. 1134–1140.
- [25] Y. Xia, H. Li, T. Qin, N. Yu, T. Y. Liu, Thompson sampling for budgeted multi-armed bandits, in: *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, 2015, pp. 3960–3966.
- [26] S. Lei, S. Maniu, L. Mo, R. Cheng, P. Senellart, Online influence maximization, in: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015, pp. 645–654.

- [27] S. Mihara, S. Tsugawa, H. Ohsaki, Influence maximization problem for unknown social networks, in: Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015, 2015, pp. 1539–1546.
- [28] S. Eshghi, S. Maghsudi, V. Restocchi, S. Stein, L. Tassiulas, Efficient Influence Maximization under Network Uncertainty, in: INFOCOM 2019 - IEEE Conference on Computer Communications Workshops, INFOCOM WKSHPs 2019, 2019, pp. 365–371.
- [29] K. Arulkumaran, M. P. Deisenroth, M. Brundage, A. A. Bharath, Deep reinforcement learning: A brief survey, IEEE Signal Processing Magazine 34 (2017) 26–38.
- [30] I. Bello, H. Pham, Q. V. Le, M. Norouzi, S. Bengio, Neural combinatorial optimization with reinforcement learning, in: Proceedings of 5th International Conference on Learning Representations, 2017, pp. 1–5.
- [31] J. J. Q. Yu, W. Yu, J. Gu, Online vehicle routing with neural combinatorial optimization and deep reinforcement learning, IEEE Transactions on Intelligent Transportation Systems 20 (2019) 3806–3817.
- [32] H. Kamarthi, P. Vijayan, B. Wilder, B. Ravindran, M. Tambe, Influence Maximization in Unknown Social Networks : Learning Policies for Effective Graph Sampling, in: Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems, 2020, pp. 575–583.
- [33] S. Wu, Q. Bai, B. H. Kang, Adaptive Incentive Allocation for Influence-Aware Proactive Recommendation, in: PRICAI 2019: Trends in Artificial Intelligence, 2019, pp. 649–661.
- [34] W. Hamilton, Z. Ying, J. Leskovec, Inductive Representation Learning on Large Graphs, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 1024–1034.
- [35] R. Ying, J. You, C. Morris, X. Ren, W. L. Hamilton, J. Leskovec, Hierarchical Graph Representation Learning with Differentiable Pooling, in: Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2018, pp. 4805–4815.
- [36] S. Fujimoto, H. V. Hoof, D. Meger, Addressing Function Approximation Error in Actor-Critic Methods, in: Proceedings of the 35th International Conference on Machine Learning, 2018, pp. 1582–1591.
- [37] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, D. Wierstra, Continuous control with deep reinforcement learning, in: Proceedings of 4th International Conference on Learning Representations, 2016, pp. 1–14.

- [38] D. Lusseau, K. Schneider, O. J. Boisseau, P. Haase, E. Slooten, S. M. Dawson, The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations, *Behavioral Ecology and Sociobiology* 54 (2003) 396–405.
- [39] J. Leskovec, J. J. McAuley, Learning to discover social circles in ego networks, in: *Proceedings of the 25th International Conference on Neural Information Processing Systems*, 2012, pp. 539–547.
- [40] J. Leskovec, D. Huttenlocher, J. Kleinberg, Signed networks in social media, in: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2010, p. 1361–1370.
- [41] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: *Proceedings of 3rd International Conference on Learning Representations*, 2015, pp. 1–15.