

Deep Graph Representation Learning and Optimization For Influence Maximization

Chen Ling, Junji Jiang, Junxiang Wang, My T. Thai, Lukas Xue, Meikang Qiu, and Liang Zhao
The 40th INTERNATIONAL CONFERENCE ON MACHINE LEARNING



Full Paper

Code



EMORY
UNIVERSITY

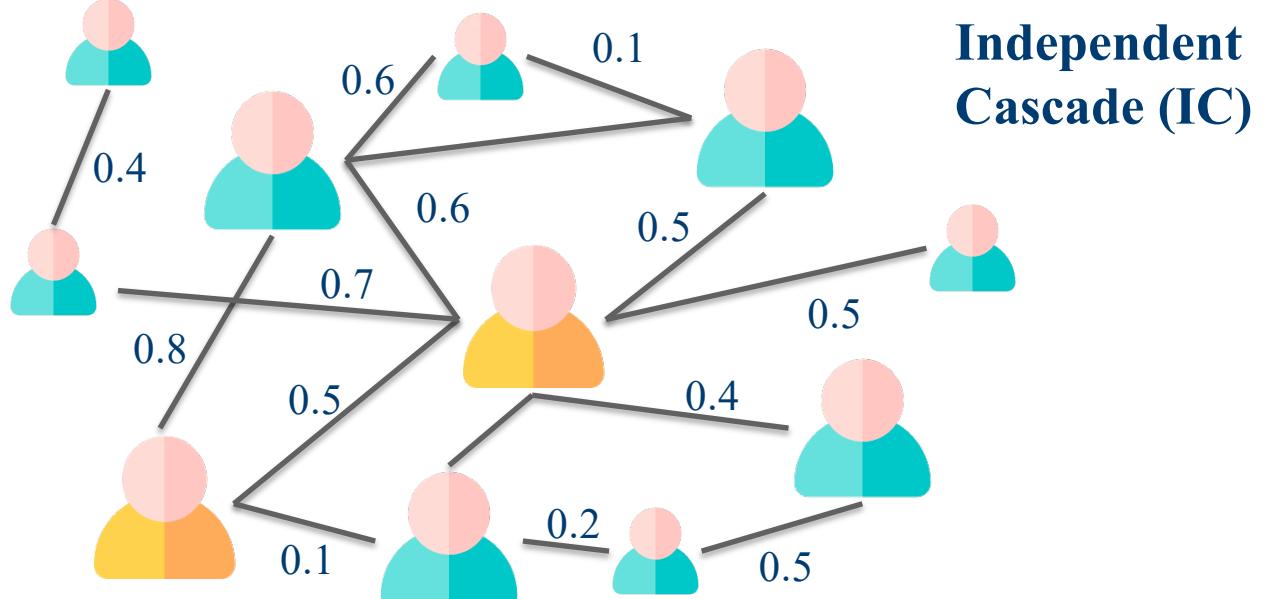
Department of
Computer Science

1. Abstract

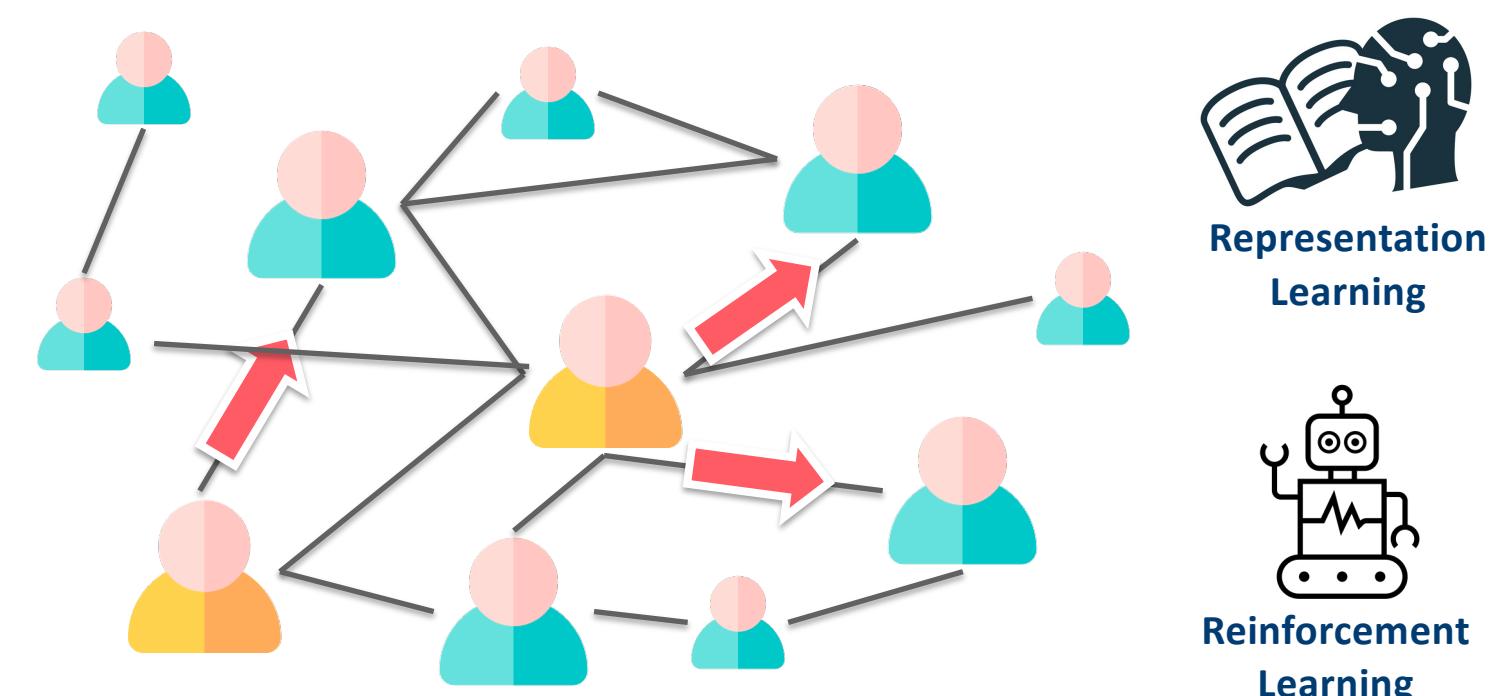
Influence maximization (IM), is formulated as selecting a set of initial users to maximize the expected number of influenced users, has seen notable improvements through traditional methods. Yet, their design and performance are near their limits. Recent learning-based IM methods provide enhanced generalization to unknown graphs but face fundamental challenges, including: 1) efficient objective function resolution; 2) representation of varied diffusion patterns; and 3) adaptation to node-centrality-constrained IM variants. To address these, we introduce a novel framework uses latent representation of seed sets to learn the diverse diffusion pattern in a data-driven, end-to-end manner. Additionally, we formulate a new objective function to infer ideal seed sets under flexible node-centrality budget constraints.

2. Challenges and Motivation

Traditional IM solutions achieve great performance with approximation ratio guarantee, but their solutions are less flexible since they need explicit diffusion model as input.



Learning-based IM methods are proposed to tackle the **NP-Hard** problem. They improve the solution generalization ability but also bring difficulties in 1) objective solving efficiency, 2) identifying diffusion process, and 3) adapting their solutions under various constraints.



3. Our Solution: DeepIM

Key Insights. In this paper, we propose a novel framework **DeepIM** to tackle the IM problem in a more robust and generalized way than existing learning-based IM methods.

1. To characterize the complex nature of the seed set, we propose to characterize the **probability** of the seed set and directly search for a more optimal seed set in continuous space.
2. We further design a learning-based diffusion model with **knowledge distillation** to characterize the underlying diffusion dynamics in an end-to-end manner.
3. we develop a generic seed set inference framework to **directly optimize and generate** set embeddings under a uniform budget constraint in order to meet various IM variant problems.

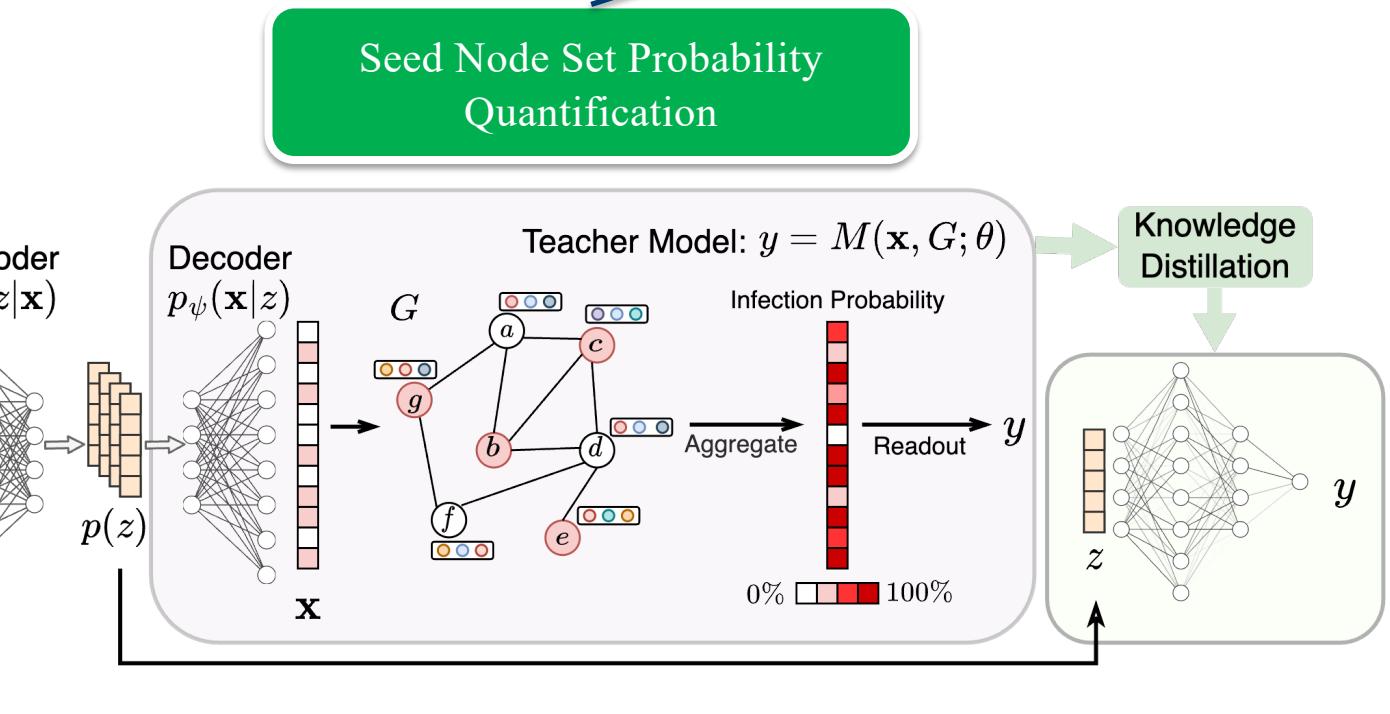
4. Training Phase

Problem Formulation. Given a graph $G = (V, E)$. The **seed node set** is defined over V as $\mathbf{x} = \{0, 1\}^{|V|}$, $x_i = 1$ or 0 denotes seed node or not. The **total number** of infected nodes is denoted as $y \in \mathbb{R}^+$. IM aims at selecting $\tilde{\mathbf{x}} = \arg \max_{|\mathbf{x}| \leq k} M(\mathbf{x}, G; \theta)$, where $y = M(\cdot)$ denotes a diffusion estimation function.

1. To build an efficient objective function:
→ Characterize the probability of the seed node set $p(\mathbf{x})$ over \mathbf{x} given the graph G .

2. To model the underlying diffusion pattern:
→ Two different diffusion models $M(\cdot)$ and $M_S(\cdot)$ to characterize the diversified diffusion dynamics with efficiency and efficacy guarantee.

$$\mathcal{L}_{\text{train}} = \min_{\theta, \lambda, \psi, \phi} \mathbb{E} \left[-\log [p_\theta(y|\mathbf{x}, G)] - \log [p_\lambda(y_s|z)] - \log [p_\psi(\mathbf{x}|z) \cdot (p_\phi(z|\mathbf{x}))] \right], \text{ s.t. } \theta \geq 0.$$



5. Inference Phase

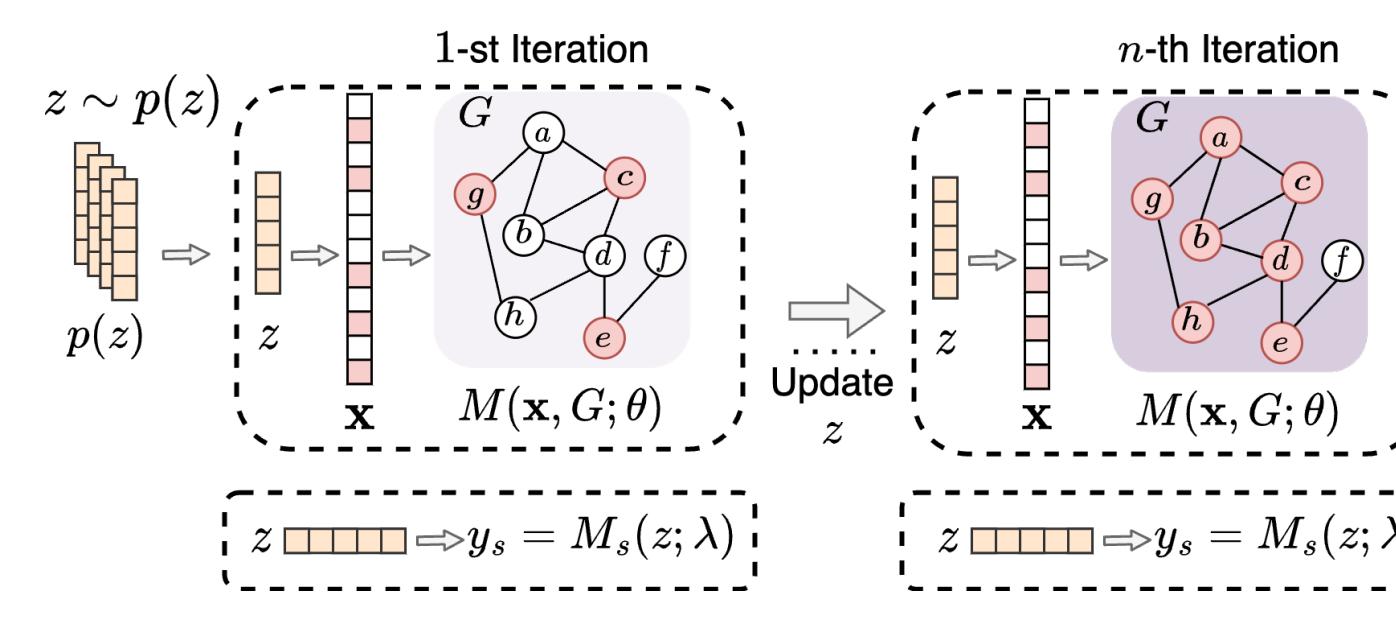
Optimization for Seed Set Inference. We propose to search the optimal seed node set $\tilde{\mathbf{x}}$ in the lower-dimensional latent space $p(z)$ by iteratively optimizing the latent variable z sampled from $p(z)$.

3. To meet different node-centrality-constraints:
→ Propose a novel objective function that can be coupled with multiple constraints for seed node set inference, which can adapt to different IM application schemes.

$$\mathcal{L}_{\text{pred}} = \max_z \mathbb{E} [p_\theta(y|\mathbf{x}, G) \cdot p_\psi(\mathbf{x}|z)], \\ \text{s.t. } \sum_{i=0}^{|V|} \mathcal{F}(v_i, G) \cdot x_i \leq k,$$

$\sum_{i=0}^{|V|} \mathcal{F}(v_i, G) \cdot x_i$ is a generalized budget constraint applied on individual nodes, and k is the actual budget (e.g., the # of nodes or the cost of selecting nodes).

$$\mathcal{L}_{\text{pred}} = \min_z \left[-\log \left[\prod_{i=0}^{|V|} f_\psi(z_i)^{x_i} (1 - f_\psi(z_i))^{1-x_i} \right] + \|\tilde{y} - y\|_2^2 \right] \text{ s.t. } \sum_{i=0}^{|V|} \mathcal{F}(v_i, G) \cdot x_i \leq k,$$



6. Experiments

Data and Baselines. We conducted experiments on six real-world networks and one synthetic network. We compare our method with both traditional and learning-based approaches.

	Digg	Weibo	Power Grid	Network Science	Cora-ML	Jazz	Synthetic
Nodes	279,613	2,251,166	4,941	1,565	2,810	198	50,000
Edges	1,170,689	225,877,808	6,594	13,532	7,981	2,742	250,000

The primary purpose is to evaluate the portion of influence spread y over the entire node set.

Methods	Cora-ML			Network Science			Power Grid			Jazz			Synthetic			Digg			Weibo									
	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%								
IMM	8.1	26.2	37.3	50.2	5.2	16.8	27.0	45.7	4.3	17.4	31.5	51.1	2.6	20.1	31.4	42.8	9.2	26.2	36.3	51.6	7.4	18.4	32.8	49.6	9.5	23.8	36.4	50.5
OPIM	13.4	26.9	37.4	50.9	6.6	19.4	28.9	45.7	5.7	17.7	29.7	50.1	2.4	20.4	34.4	46.8	9.6	25.3	36.6	51.7	7.6	18.5	32.9	48.9	9.7	23.7	36.6	50.3
SubSIM	10.1	25.7	36.8	51.1	4.8	15.4	27.5	44.8	4.6	19.2	31.7	50.2	3.6	18.8	37.6	44.7	9.5	26.7	36.5	51.3	7.5	18.9	33.3	49.4	9.3	23.1	36.5	50.6
OIM	8.9	27.6	38.0	51.3	4.2	16.7	26.5	48.2	5.7	17.5	31.9	50.8	2.0	18.5	36.3	42.2	9.6	26.2	36.1	51.5	7.8	18.2	33.1	49.6	-	-	-	-
IMINFECTOR	13.6	26.8	37.7	50.4	5.4	17.9	27.8	47.6	5.4	18.2	31.0	50.9	3.6	19.7	32.7	45.9	9.1	26.4	36.1	51.5	7.9	18.6	33.5	49.8	9.4	23.5	36.9	50.3
PIANO	13.6	26.8	37.7	50.4	5.4	17.9	27.8	47.6	5.4	18.2	31.0	50.9	3.6	19.7	32.7	45.9	9.1	26.4	36.1	51.5	7.9	18.6	33.5	49.8	9.4	23.5	36.9	50.3
ToupleGDD	10.6	27.5	38.5	51.3	6.3	17.8	28.3	50.5	5.4	19.3	31.6	51.3	3.3	20.4	37.2	45.7	9.5	26.8	37.1	51.4	7.8	18.8	33.7	49.8	9.4	23.5	36.9	50.8
DeepIM _s	14.4	28.1	39.6	52.4	7.8	20.9	31.5	51.2	6.3	21.0	32.5	52.4	4.9	23.3	41.5	49.9	11.6	27.4	38.7	52.1	8.4	19.3	34.2	51.3	11.2	26.5	37.9	51.8
DeepIM	14.4	28.1	39.6	52.4	7.8	20.9	31.5	51.2	6.3	21.0	32.5	52.4	4.9	23.3	41.5	49.9	11.6	27.4	38.7	52.1	8.4	19.3	34.2	51.3	11.2	26.5	37.9	51.8

Table 2. Performance comparison under IC diffusion pattern. — indicates out-of-memory error. (Best is highlighted with bold.)

Methods	Cora-ML			Network Science			Power Grid			Jazz			Synthetic			Digg			Weibo									
	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%	1%	5%	10%	20%								
IMM	1.7	34.8	52.2	66.4	2.5	11.9	18.1	33.6	4.6	19.9	31.7	56.9	1.4	5.7	13.4	24.5	1.1	5.2	13.1	66.9	2.4	10.8	34.5	55.6	1.6	6.7	19.3	45.2
OPIM	2.3	36.9	51.2	71.5	1.6	12.0	18.2	34.4																				