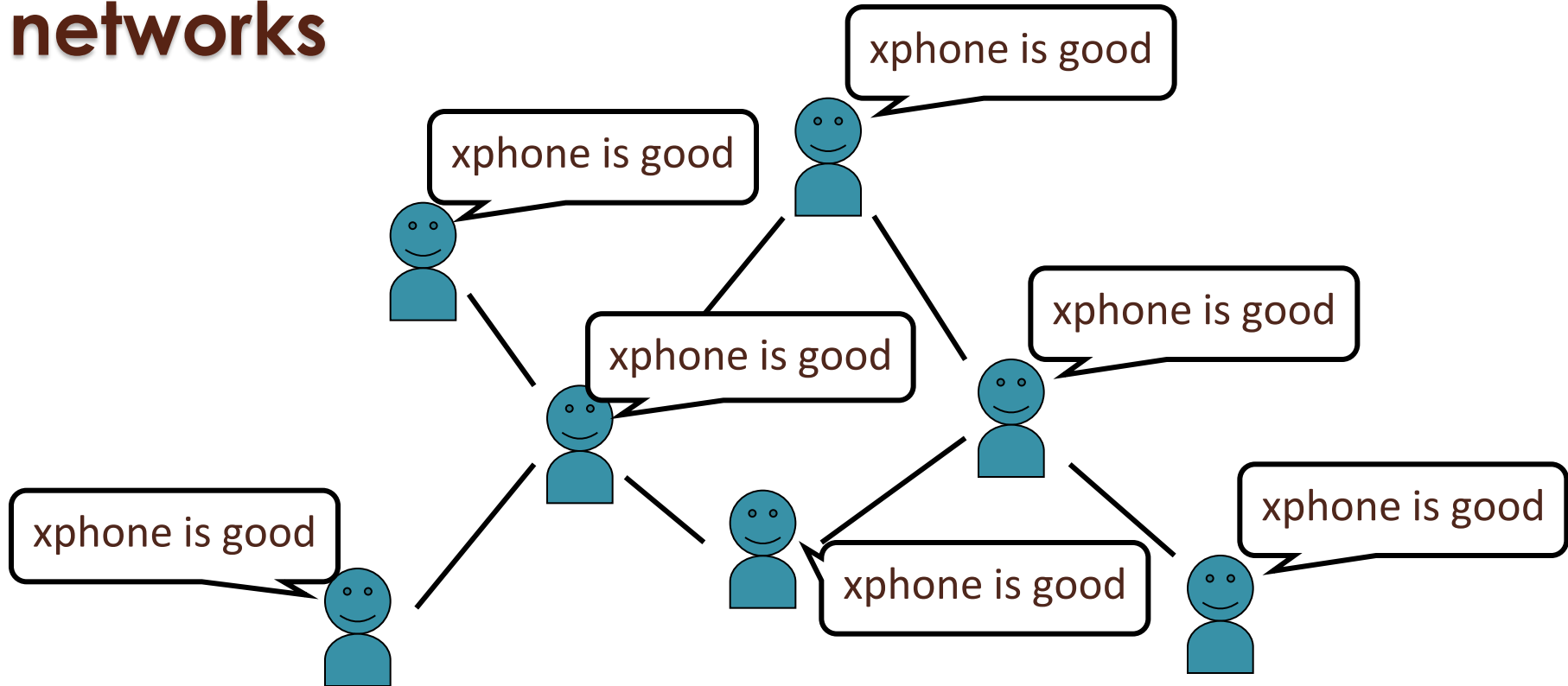


Part I → Part II → Part III → Part IV → Part V

Influence Maximization



Word-of-mouth (WoM) effect in social networks



- Word-of-mouth (viral) marketing is believed to be a promising marketing strategy.
- Increasing popularity of online social networks may enable large scale viral marketing

Outline

- Diffusion/propagation models and the Influence Maximization (IM) problem
- The theory: greedy methods to optimize submodular functions
- Scalable influence maximization



Diffusion/Propagation Models and the Influence Maximization (IM) Problem



The first definition of IM problem: A Markov random fields formulation

- Each node i has random variable X_i , indicating bought the product or not,
 $\mathbf{X} = \{X_1, \dots, X_n\}$
- Markov random field formation: X_i depends on its neighbors' actions $\mathbf{N}_i \subseteq \mathbf{X}$
- Marketing action $\mathbf{M} = \{M_1, \dots, M_n\}$
- **Problem**: find a choice of \mathbf{M} that maximizes the revenue obtained from the result of \mathbf{X}

Discrete diffusion models

- Static social network: $G = (V, E)$
 - V : set of nodes, representing individuals
 - E : set of edges, representing directed influence relationships
- Influence diffusion process
 - **Seed set S** : initial set of nodes selected to start the diffusion
 - **Node activations**: Nodes are activated starting from the seed nodes, in discrete steps and following certain stochastic models
 - **Influence spread $\sigma(S)$** : expected number of activated nodes when the diffusion process starting from the seed set S ends

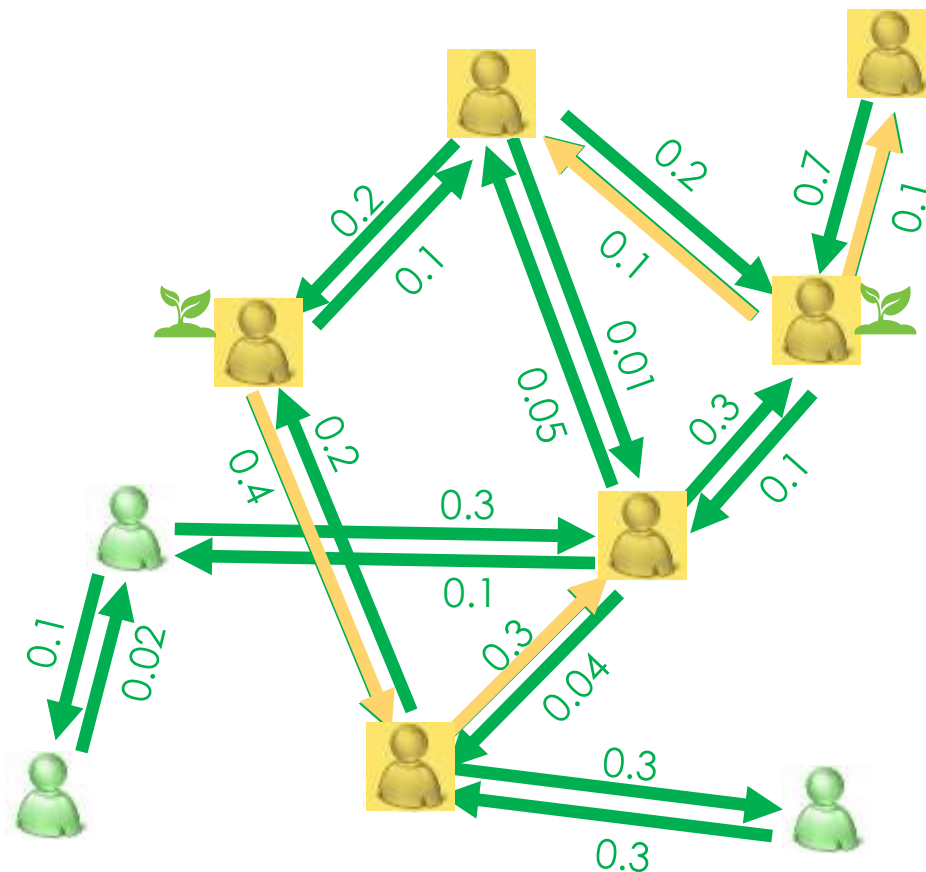
Major stochastic diffusion models

- **Independent cascade** (IC) model
- **Linear threshold** (LT) model
- General threshold model
- Others
 - Voter model
 - Heat diffusion model



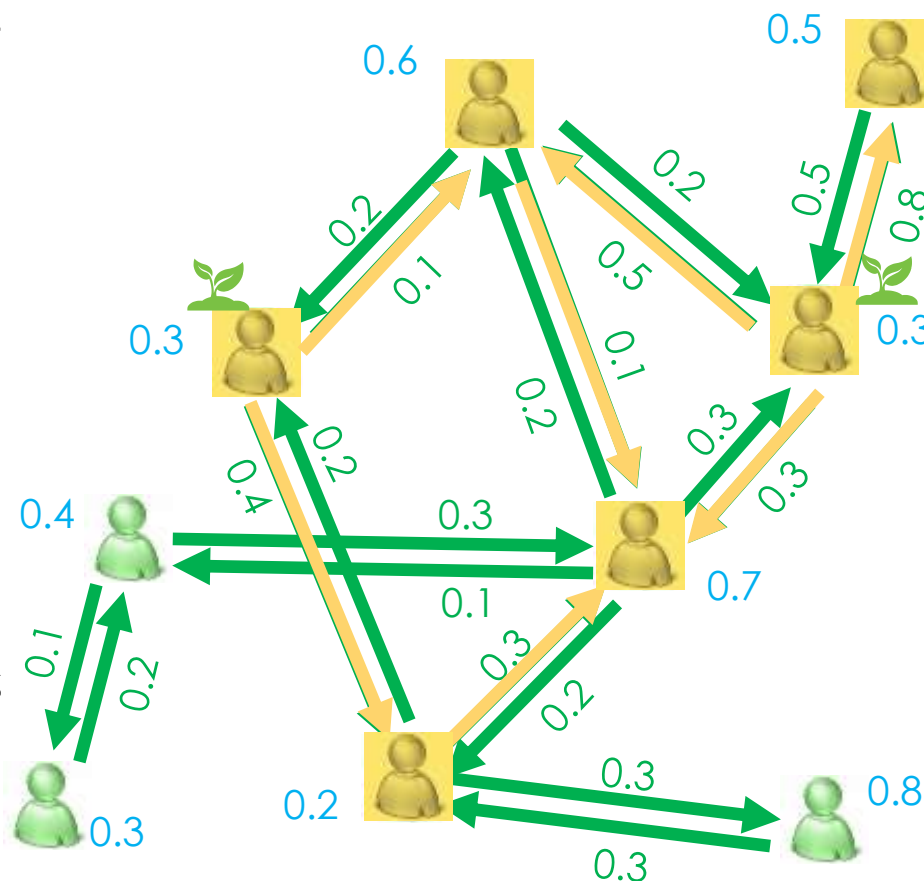
Independent cascade model

- Each edge (u, v) has a propagation probability $p(u, v)$
- Initially some seed nodes are activated
- At each step t , each node u activated at step $t - 1$ activates its neighbor v with probability $p(u, v)$
- once activated, stay activated



Linear threshold model

- Each edge (u, v) has weight $w(u, v)$:
 - when $(u, v) \notin E, w(u, v) = 0$
 - $\sum_u w(u, v) \leq 1$
- Each node v selects a threshold $\theta_v \in [0, 1]$ uniformly at random
- Initially some seed nodes are activated
- At each step, node v checks if the weighted sum of its active neighbors is greater than its threshold θ_v , if so v is activated
- once activated, stay activated



Influence maximization

- Given a social network, a diffusion model with given parameters, and a number k , find a seed set S of at most k nodes such that the influence spread of S is maximized.
- May be further generalized:
 - Instead of k , given a budget constraint and each node has a cost of being selected as a seed
 - Instead of maximizing influence spread, maximizing a (submodular) function of the set of activated nodes

Key takeaways for diffusion models and IM definition

- stochastic diffusion models
 - IC model reflects simple contagion
 - LT model reflects complex contagion (activation needs social affirmation from multiple sources [Centola and Macy, AJS 2007])
- maximization objective focuses on expected influence spread
 - others to be considered later



Time for **Theory:**

Optimizing Submodular Functions



[Credit: Rico3244 @ DeviantArt](#)



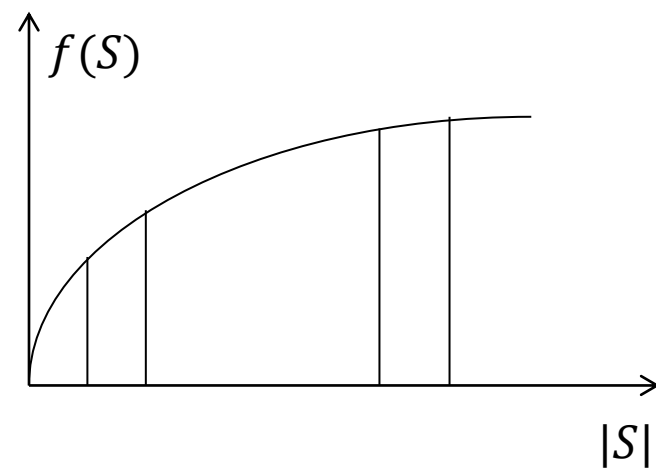
Hardness of influence maximization

- Influence maximization under both IC and LT models are NP hard
 - IC model: reduced from k-max cover problem
 - LT model: reduced from vertex cover problem
- Need approximation algorithms



Optimizing submodular functions

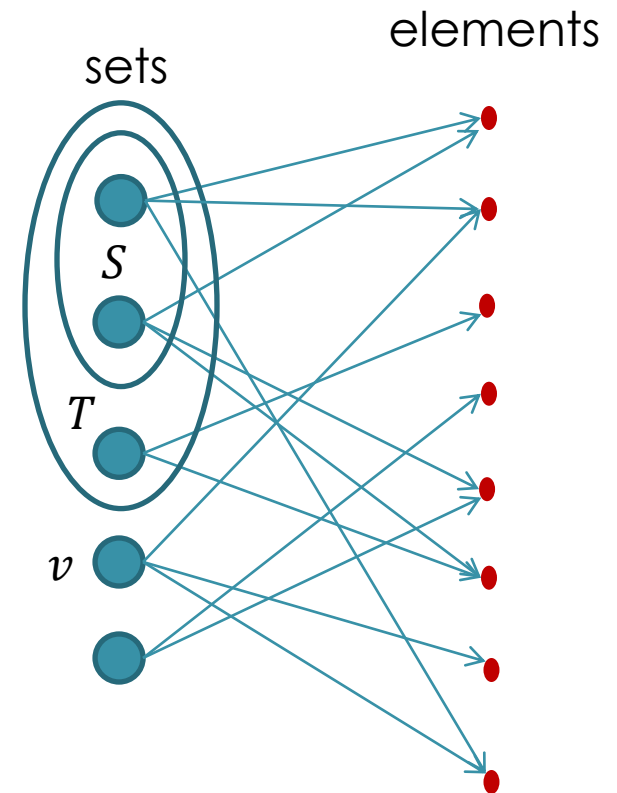
- **Sumodularity** of set functions $f: 2^V \rightarrow R$
 - for all $S \subseteq T \subseteq V$, all $v \in V \setminus T$,
 $f(S \cup \{v\}) - f(S) \geq f(T \cup \{v\}) - f(T)$
 - diminishing marginal return
 - an equivalent form: for all $S, T \subseteq V$
 $f(S \cup T) + f(S \cap T) \leq f(S) + f(T)$



- **Monotonicity** of set functions f : for all $S \subseteq T \subseteq V$,
 $f(S) \leq f(T)$

Example of a submodular function and its maximization problem

- set coverage
 - each entry u is a subset of some base elements
 - coverage $f(S) = |\bigcup_{u \in S} u|$
 - $f(S \cup \{v\}) - f(S)$: additional coverage of v on top of S
- k -max cover problem
 - find k subsets that maximizes their total coverage
 - NP-hard
 - special case of IM problem in IC model



Greedy algorithm for submodular function maximization

- 1: initialize $S = \emptyset$;
- 2: for $i = 1$ to k do
- 3: select
 $u = \operatorname{argmax}_{w \in V \setminus S} [f(S \cup \{w\}) - f(S)]$
- 4: $S = S \cup \{u\}$
- 5: end for
- 6: output S

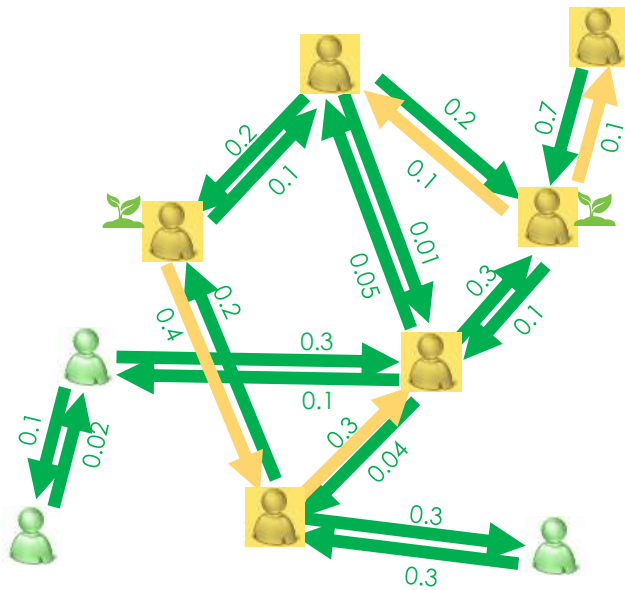


Property of the greedy algorithm

- Theorem: If the set function f is monotone and submodular with $f(\emptyset) = 0$, then the greedy algorithm achieves $(1 - 1/e)$ approximation ratio, that is, the solution S found by the greedy algorithm satisfies:
 - $f(S) \geq \left(1 - \frac{1}{e}\right) \max_{S' \subseteq V, |S'|=k} f(S')$

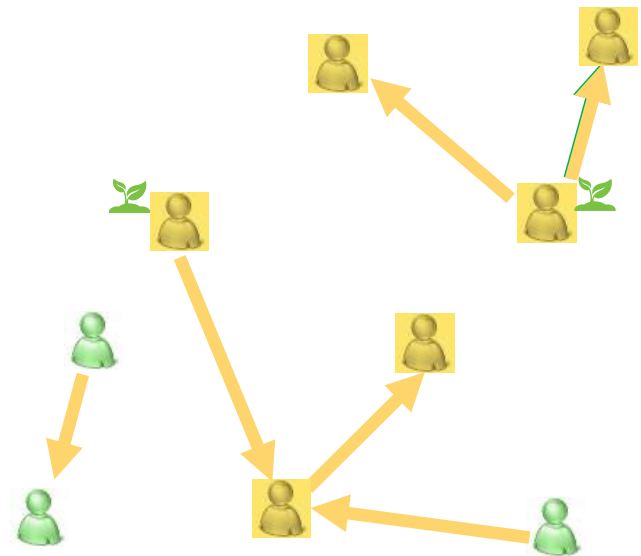
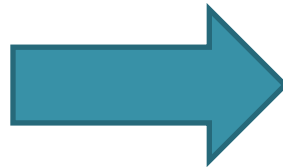
Submodularity of influence diffusion models

- Based on equivalent live-edge graphs



diffusion dynamic

$\Pr(\text{set A is activated given seed set S})$



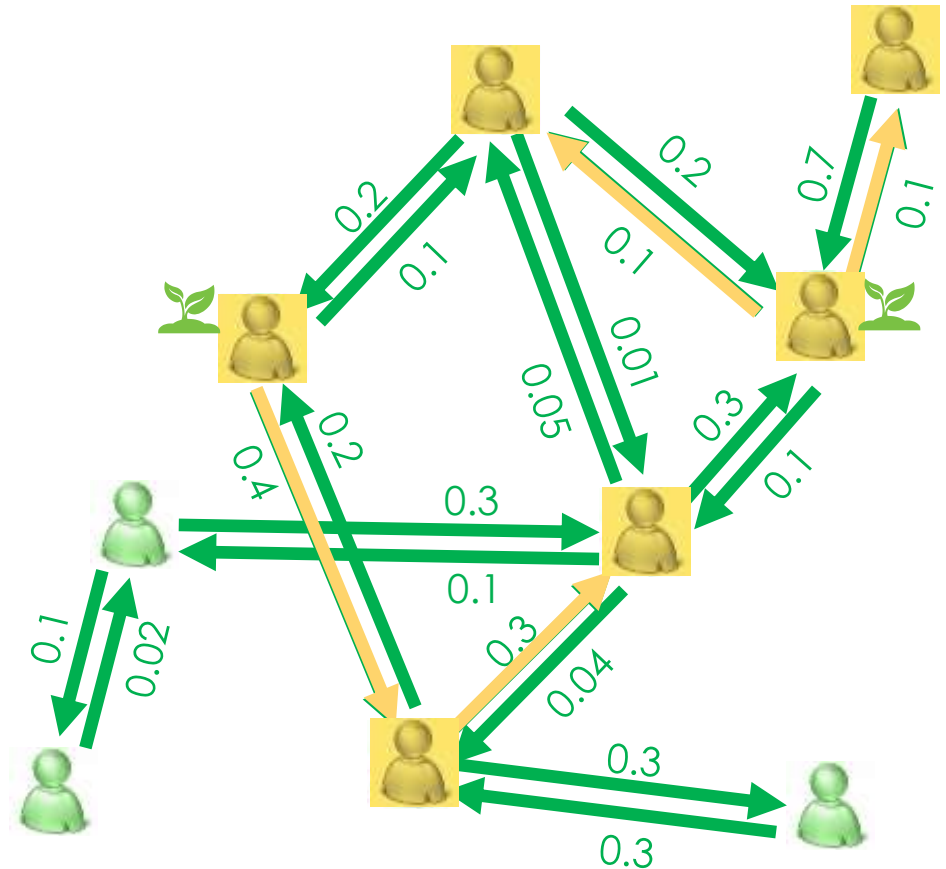
random live-edge graph:
edges are randomly removed

$\Pr(\text{set A is reachable from S in random live-edge graph})$



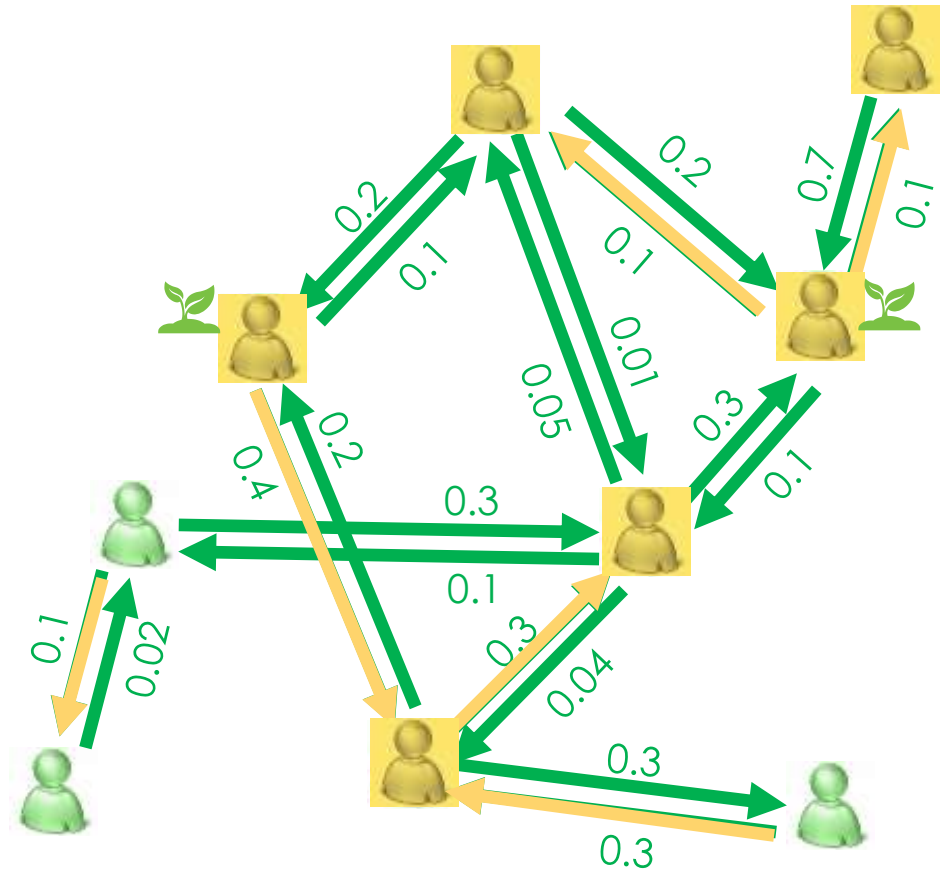
(Recall) active node set via IC diffusion process

- yellow node set is the active node set after the diffusion process in the independent cascade model



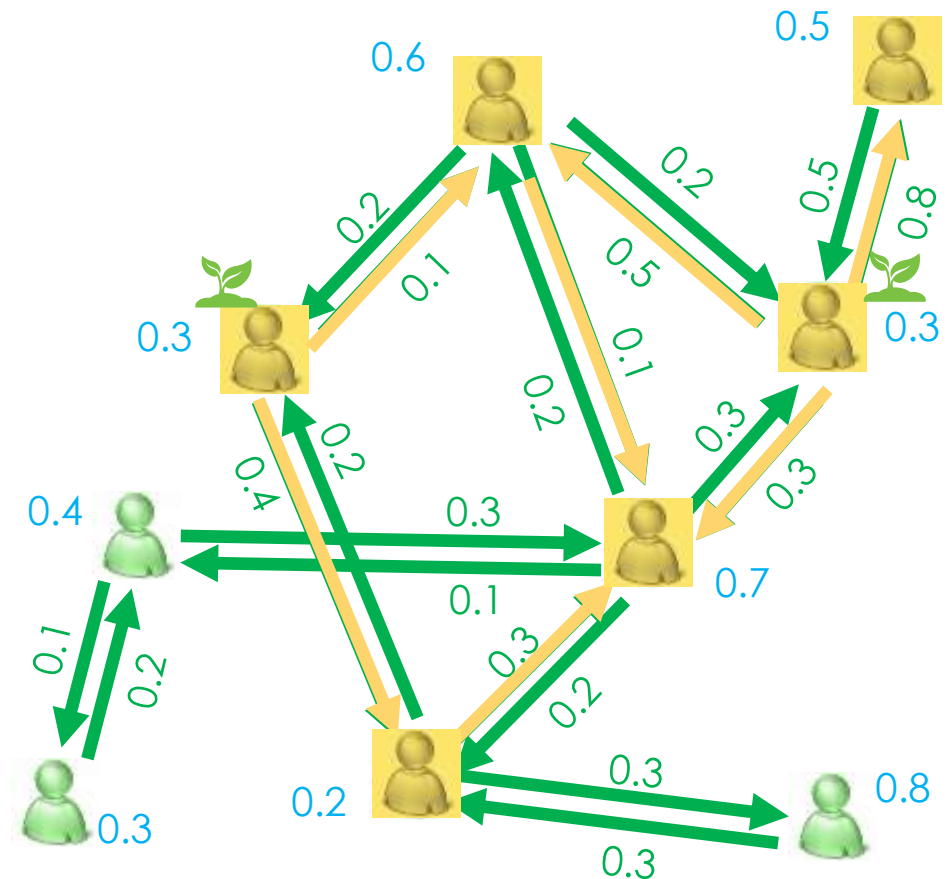
Random live-edge graph for the IC model and its reachable node set

- random live-edge graph in the IC model
 - each edge is independently selected as live with its propagation probability
- yellow node set is the active node set reachable from the seed set in a random live-edge graph
- Equivalence is straightforward



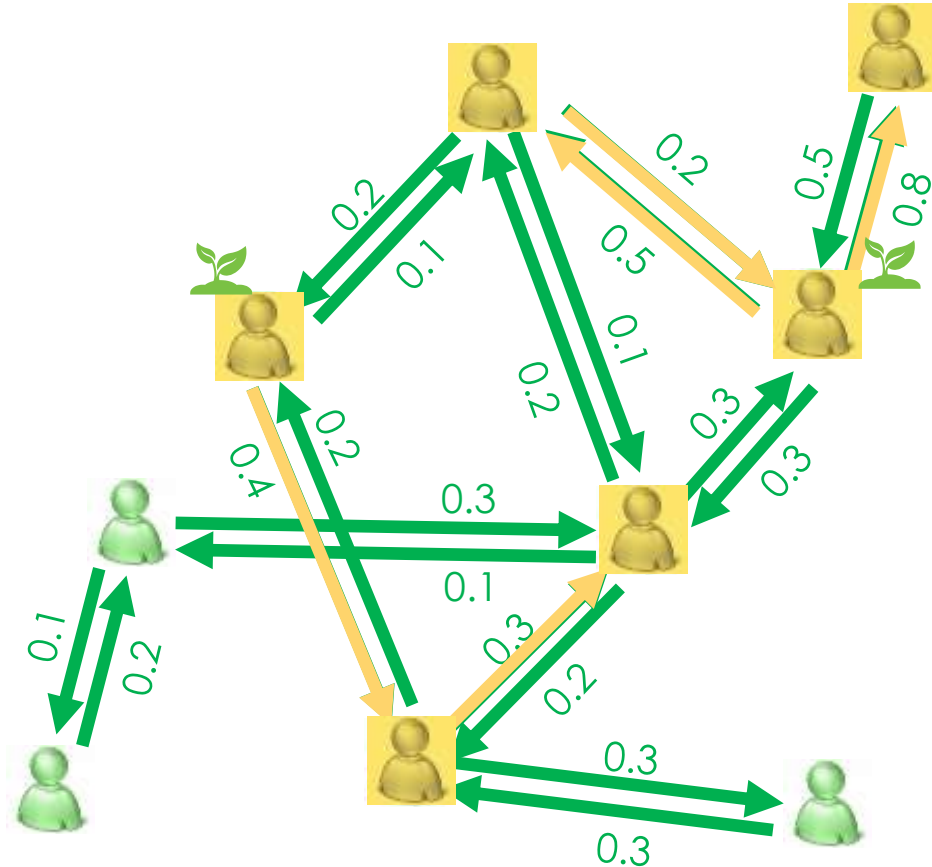
(Recall) active node set via LT diffusion process

- yellow node set is the active node set after the diffusion process in the linear threshold model



Random live-edge graph for the LT model and its reachable node set

- random live-edge graph in the LT model
 - each node select at most one incoming edge, with probability proportional to its weight
- yellow node set is the active node set reachable from the seed set in a random live-edge graph
- equivalence is based on uniform threshold selection from $[0,1]$, and linear weight addition

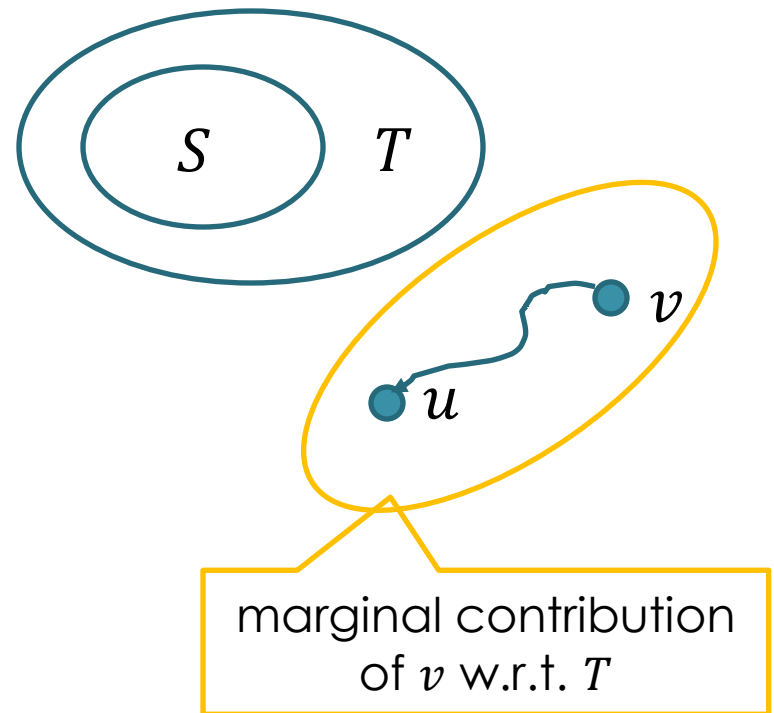


Submodularity of influence diffusion models (cont'd)

- Influence spread of seed set S , $\sigma(S)$:
$$\sigma(S) = \sum_{G_L} \Pr(G_L) |R(S, G_L)|,$$
 - G_L : a random live-edge graph
 - $\Pr(G_L)$: probability of G_L being generated
 - $R(S, G_L)$: set of nodes reachable from S in G_L
- To prove that $\sigma(S)$ is submodular, only need to show that $|R(\cdot, G_L)|$ is submodular for any G_L
 - submodularity is maintained through linear combinations with non-negative coefficients

Submodularity of influence diffusion models (cont'd)

- Submodularity of $|R(\cdot, G_L)|$
 - for any $S \subseteq T \subseteq V$,
 $v \in V \setminus T$,
 - if u is reachable from v but not from T , then
 - u is reachable from v but not from S
 - Hence, $|R(\cdot, G_L)|$ is submodular
- Therefore, influence spread $\sigma(S)$ is submodular in both IC and LT models



General threshold model

- Each node v has a threshold function
$$f_v: 2^V \rightarrow [0,1]$$
- Each node v selects a threshold $\theta_v \in [0,1]$ uniformly at random
- If the set of active nodes at the end of step $t - 1$ is S , and $f_v(S) \geq \theta_v$, v is activated at step t
- reward function $r(A(S))$: if $A(S)$ is the final set of active nodes given seed set S , $r(A(S))$ is the reward from this set
- generalized influence spread:

$$\sigma(S) = E[r(A(S))]$$

[Kempe, Kleinberg and Tardos, KDD 2003]

IC and LT as special cases of general threshold model

- LT model
 - $f_v(S) = \sum_{u \in S} w(u, v)$
 - $r(S) = |S|$
- IC model
 - $f_v(S) = 1 - \prod_{u \in S} (1 - p(u, v))$
 - $r(S) = |S|$



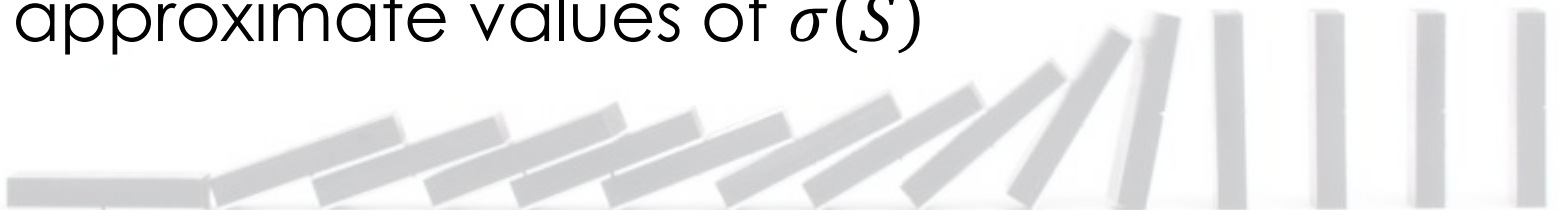
Submodularity in the general threshold model

- Theorem [Mossel & Roch STOC 2007]:
 - In the general threshold model,
 - if for every $v \in V$, $f_v(\cdot)$ is monotone and submodular with $f_v(\emptyset) = 0$,
 - and the reward function $r(\cdot)$ is monotone and submodular,
 - then the general influence spread function $\sigma(\cdot)$ is monotone and submodular.
- Local submodularity implies global submodularity



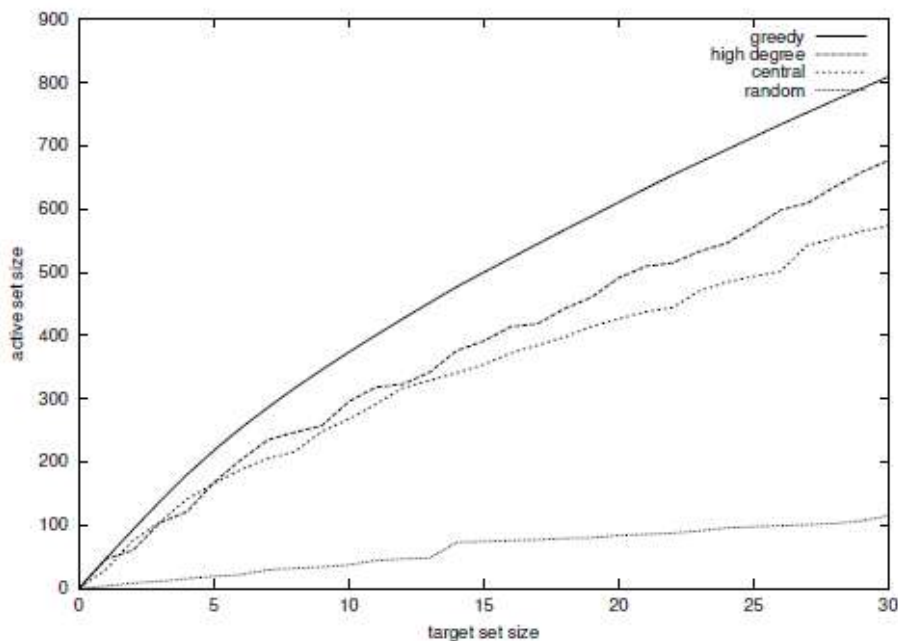
The greedy approximation algorithm for the IM problem

- $\left(1 - \frac{1}{e} - \varepsilon\right)$ -approximation greedy algorithm
 - Use general greedy algorithm framework
 - However, need evaluation of influence spread $\sigma(S)$, which is shown to be #P-hard
 - Use Monte Carlo simulations to estimate $\sigma(S)$ to arbitrary accuracy with high probability
 - For any $\varepsilon > 0$, there exists $\gamma > 0$, s.t. $\left(1 - \frac{1}{e} - \varepsilon\right)$ -approximation can be achieved with $1 - \gamma$ approximate values of $\sigma(S)$

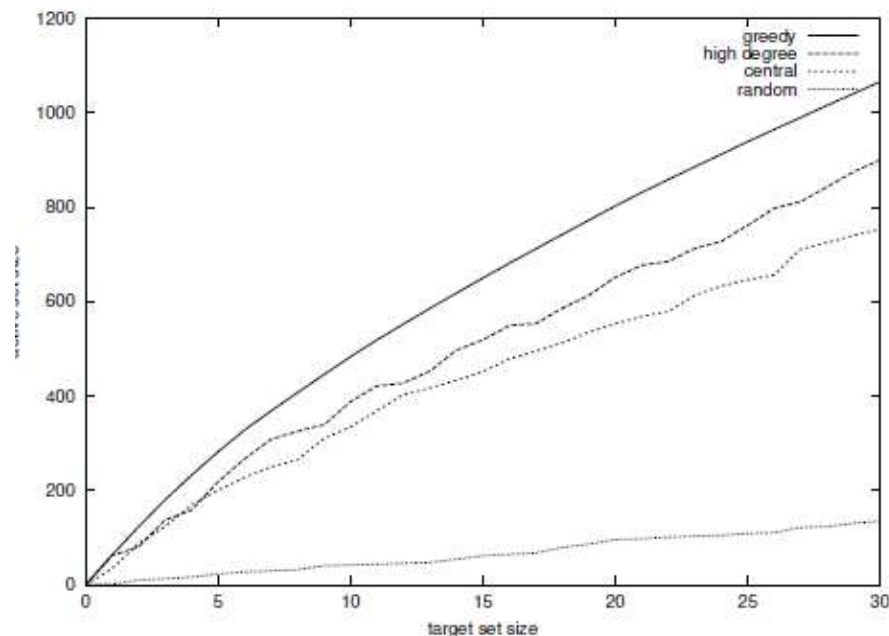


Performance of greedy algorithm

weighted cascade model



linear threshold model



- coauthorship graph from arxiv.org, high energy physics section, $n=10748$, $m \approx 53000$
- allow parallel edges, $c_{u,v}$ = number of edges between u and v
- weighted cascade: $p(u, v) = 1 - \left(1 - \frac{1}{d_v}\right)^{c_{u,v}}$
- linear threshold: $w(u, v) = c_{u,v}/d_v$

Key takeaways for theory support

- submodularity of diffusion models
 - diminishing return
 - shared by many models, but some model extensions may not be submodular (see Part IV)
- submodular function maximization
 - greedy hill-climbing algorithm
 - approximation ratio $(1 - 1/e)$

Scalable Influence Maximization



Theory versus Practice



...the trouble about arguments is, they ain't nothing but theories, after all, and theories don't prove nothing, they only give you a place to rest on, a spell, when you are tuckered out butting around and around trying to find out something there ain't no way to find out...

There's another trouble about theories: there's always a hole in them somewheres, sure, if you look close enough.

- "Tom Sawyer Abroad", Mark Twain

Inefficiency of greedy algorithm

- It takes days to find 50 seeds for a 30K node graph
- Too many evaluations of influence spread
 - Given a graph of n nodes, each round needs evaluation of n influence spreads, totally $O(nk)$ evaluations
- Each evaluation of influence spread is hard
 - Exact computation is #P-hard, for both IC and LT models [Chen et al. KDD/ICDM 2010]
 - Monte-Carlo simulation is very slow

Scalable Influence Maximization

- Reduction on the number of influence spread evaluations.
- Batch computation of influence spread
- Scalable heuristics for influence spread computation



Lazy forward optimization

- Exploiting submodularity, significantly reduce # influence spread evaluations
 - S_{t-1} – seed set selected after round $t - 1$
 - v_t – selected as seed in round t : $S_t = S_{t-1} \cup \{v_t\}$
 - u is not a seed yet, u 's marginal gain $MG(u|S_t) = \sigma(S_t \cup \{u\}) - \sigma(S_t)$
 - by submodularity, $MG(u|S_t) \leq MG(u|S_{t-1})$
 - This implies, if $MG(u|S_{t-1}) \leq MG(v|S_t)$ for some node v , then no need to evaluate $MG(u|S_t)$ in round $t + 1$.
 - Can be implemented efficiently using max-heap
 - take the top of the heap, if it has MG for the current round, then it is the new seed;
 - else compute its MG , and re-heapify
- Often, top element after round $k - 1$ remains top in round k .
- Up to 700 X speedup

CELF++

- With each heap node u , further maintain
 - $u.mg1 = MG(u|S)$, S = current seed set.
 - $u.prev_best$ = node with max. MG in the current iteration, among nodes seen before u ;
 - $u.mg2 = MG(u|S \cup \{u.prev_best\})$;
 - $u.flag$ = iteration where $u.mg1$ was last updated;
- $u.prev_best$ chosen in current iteration \rightarrow no need to compute $MG(u|S \cup \{u.prev_best\})$ in next iteration.
- $MG(u|S \cup \{u.prev_best\})$ can be computed efficiently along with $MG(u|S)$ in one run of MC.

CELF++ (contd.)

- Pick the top node u in the heap
- $u.flag = |S| \rightarrow u.mg1$ is correct MG , thus u is the best pick of current itn; add u to S ; remove from Q ;
- else, if $u.prev_best = last_seed$ & $u.flag = |S| - 1 \rightarrow u.mg1 \leftarrow u.mg2$ (no need to compute $MG(u|S + last_seed)$)

Dataset	Running time (min)			Avg. # node lookups		
	CELF	CELF++	Gain	CELF	CELF++	Gain
Hept WC	245	178	27%	18.7	13.6	27.2%
Hept IC	5269	2082	60.5%	190.5	113.0	40.7%
Phy WC	1242	1028	17.2%	18.6	17.9	3.8%

Table 1: Comparison between CELF and CELF++. Number of seeds = 100 in all test cases.

Batch computation of influence spread

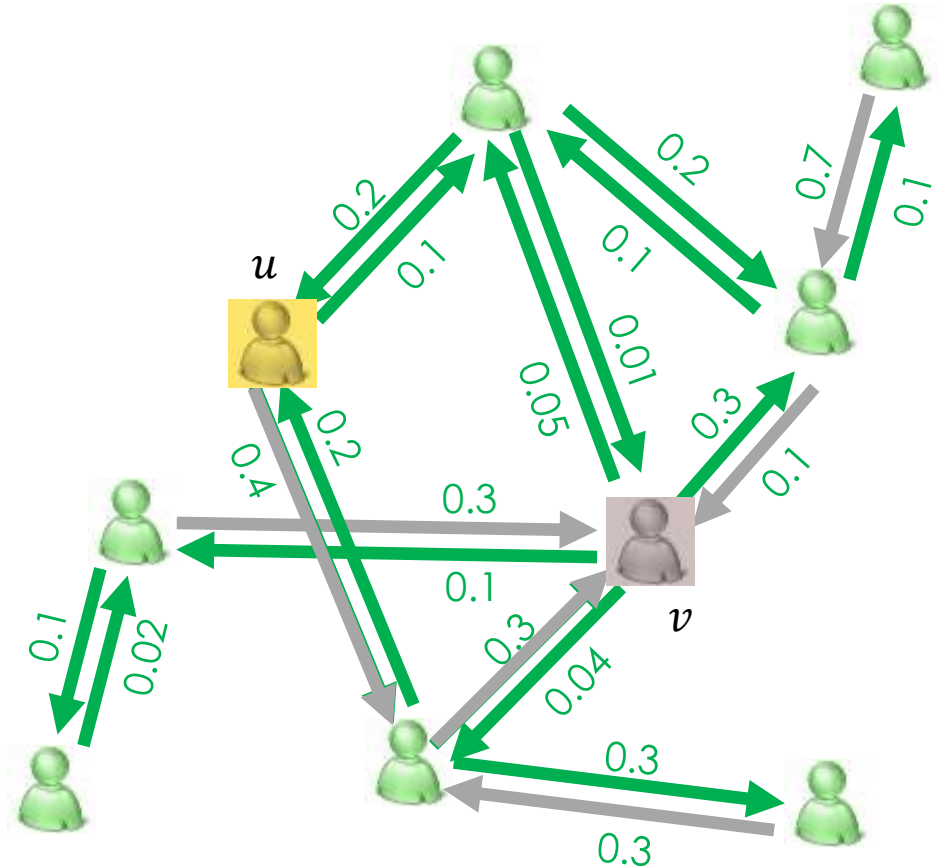
- Based on equivalence to reachable set in live-edge graphs
 - Generate a random live-edge graph
 - batch computation of size of reachable set of all nodes [Cohen, JCSS 1997]
 - repeat above to increase accuracy
- In conflict with lazy-forward optimization
 - MixedGreedy: In the first round uses reachability-based batch computation; in remaining rounds, use lazy forward optimization

Scalable heuristics for influence spread computation

- MIA: for IC model [Chen, Wang and Wang, KDD 2010]
- LDAG: for LT model [Chen, Yuan and Zhang, ICDM 2010]
- Features:
 - Focus on a local region
 - Find a graph spanner allow efficient computation
 - Batch update of marginal influence spread

Maximum Influence Arborescence (MIA) Heuristic

- Local influence regions of node v
 - For every nodes u , find the maximum influence path (MIP) from u to v , ignore it if $\Pr(path) < \lambda$ (λ is a small threshold value)
 - all MIPs to v form its maximum influence in-arborescence (MIIA)
 - MIIA can be efficiently computed
 - influence to v is computed over MIIA



MIA Heuristic II: Computing Influence through the MIA structure

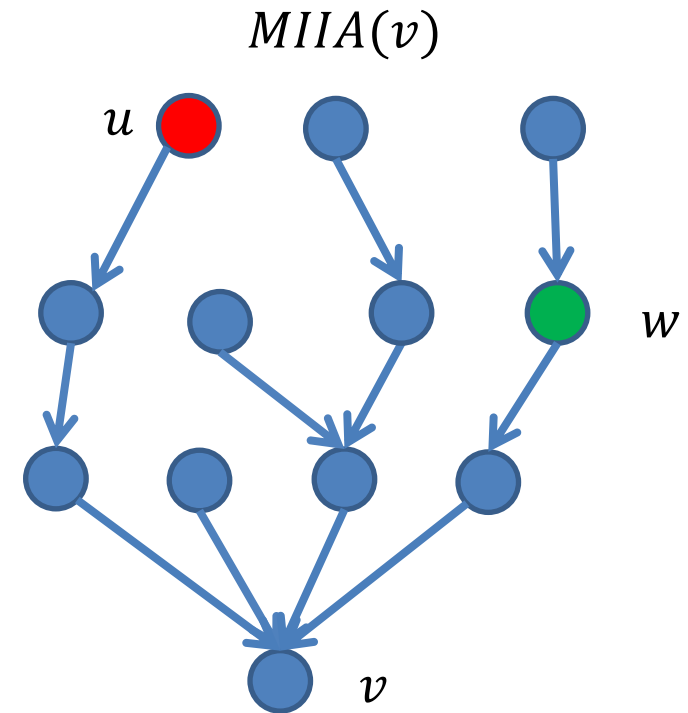
- Recursive computation of activation probability $ap(u)$ of a node u in its in-arborescence, given a seed set S ($N^{in}(u)$ is the in-neighbors of u in its in-arborescence)

Algorithm 2 Computing activation probability of u , $ap(u, S)$

```
1: if  $u \in S$  then
2:    $ap(u) = 1$ 
3: else if  $N^{in}(u) = \emptyset$  then
4:    $ap(u) = 0$ 
5: else
6:    $ap(u) = 1 - \prod_{w \in N^{in}(u)} (1 - ap(w) \cdot p(w, u))$ 
7: end if
```

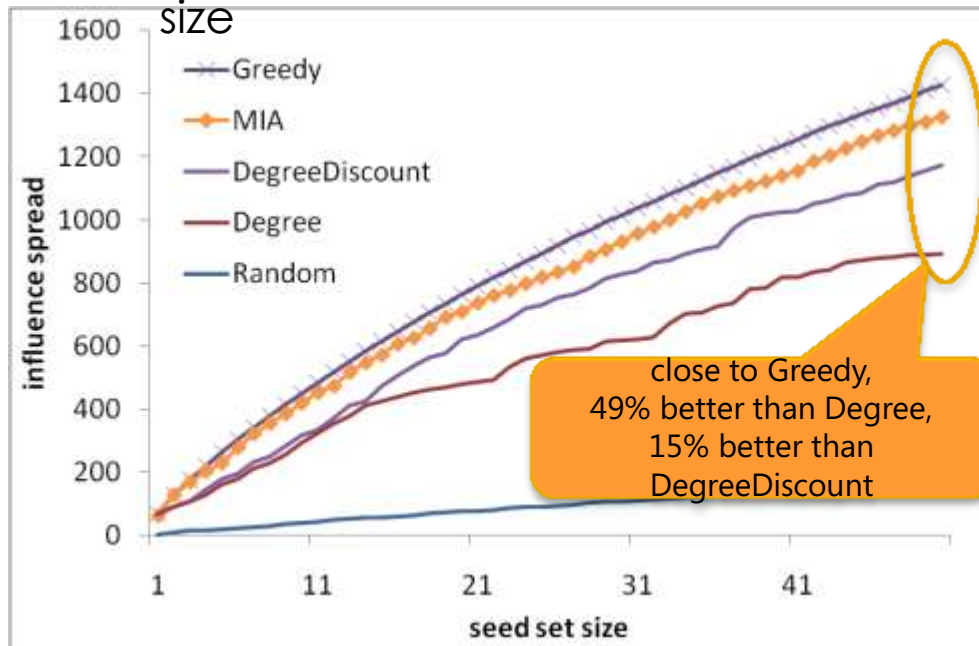
MIA Heuristic IV: Efficient updates on incremental activation probabilities

- u is the new seed in $MIIA(v)$
- Naive update: for each candidate w , redo the computation in the previous page to compute w 's marginal influence to v
 - $O(|MIIA(v)|^2)$
- Fast update: based on linear relationship of activation probabilities between any node w and root v , update marginal influence of all w 's to v in two passes
 - $O(|MIIA(v)|)$

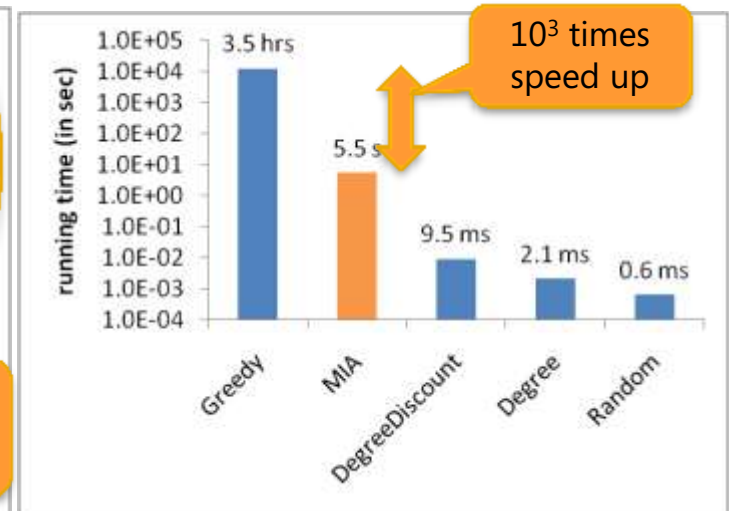


Experiment Results on MIA heuristic

Influence spread vs. seed set size



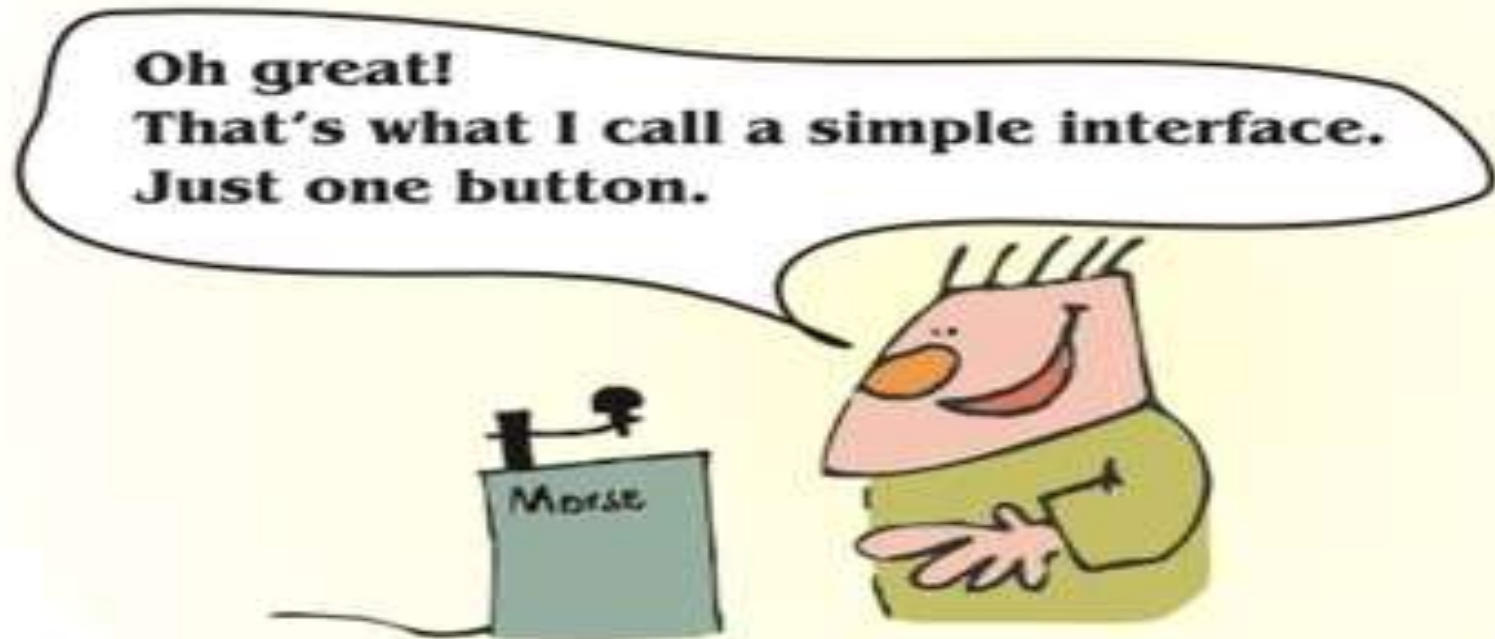
running time



Experiment setup:

- 35k nodes from coauthorship graph in physics archive
- influence probability to a node $v = 1 / (\# \text{ of neighbors of } v)$
- running time is for selecting 50 seeds

Time for some simplicity



The SimPath Algorithm

Vertex Cover Optimization

Improves the efficiency in the first iteration

Look ahead optimization

Improves the efficiency in the subsequent iterations

In lazy forward manner, in each iteration, add to the seed set, the node providing the maximum marginal gain in spread.

[Goyal, Lu, & L. ICDM 2011]

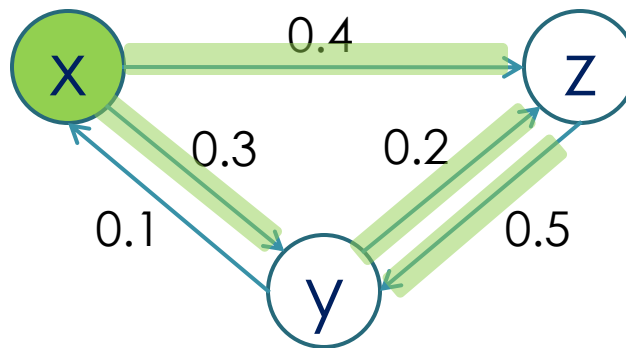


Simpath-Spread

Compute marginal gain by enumerating simple paths

Estimating Spread in SimPath (2)

- Thus, the spread of a node can be computed by enumerating simple paths starting from the node.



Influence of x on z
Influence of x on y
Influence of x on x itself
Total influence of node x is 1.96

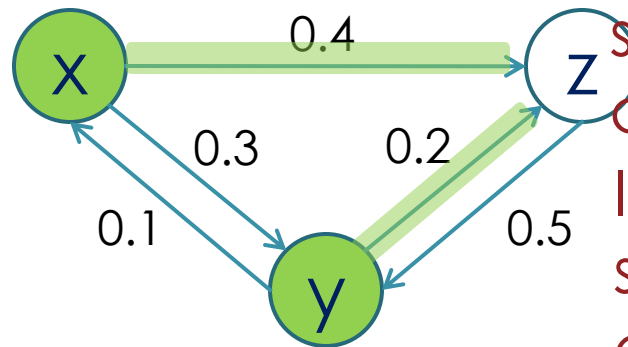
$$\begin{aligned} \text{Influence Spread of node } x &= \Upsilon_{x,x} + \Upsilon_{x,y} + \Upsilon_{x,z} \\ &= 1 + (0.3 + 0.4 * 0.5) + (0.4 + 0.3 * 0.2) = 1.96 \end{aligned}$$

Estimating Spread in SimPath (3)

Theorem 1. In the LT model, the spread of a set S is the sum of the spread of each node $u \in S$ on subgraphs induced by $V - S + u$. That is,

$$\sigma(S) = \sum_{u \in S} \sigma^{V-S+u}(u)$$

Total influence of the seed set $\{x, y\}$ is 2.6



Influence of node y in a subgraph that does not contain x

Influence of node x in a subgraph that does not contain y

Let the seed set $S = \{x, y\}$, then influence spread of S is

$$\sigma(S) = \sigma^{V-y}(x) + \sigma^{V-x}(y) = 1 + 0.4 + 1 + 0.2 = 2.6$$

Estimating Spread in SimPath (4)



Enumerating all simple paths is #P hard

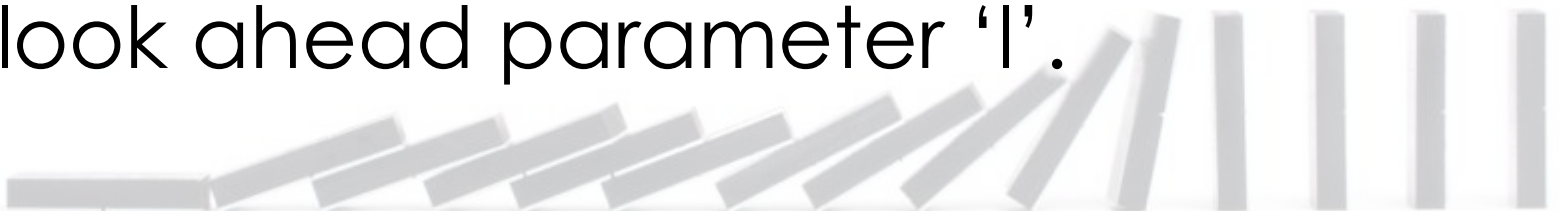
Thus, influence can be estimated
by enumerating all simple paths
starting from the seed set.

On slightly different subgraphs

The majority of influence flows in a small neighborhood.

Look Ahead Optimization (1/2)

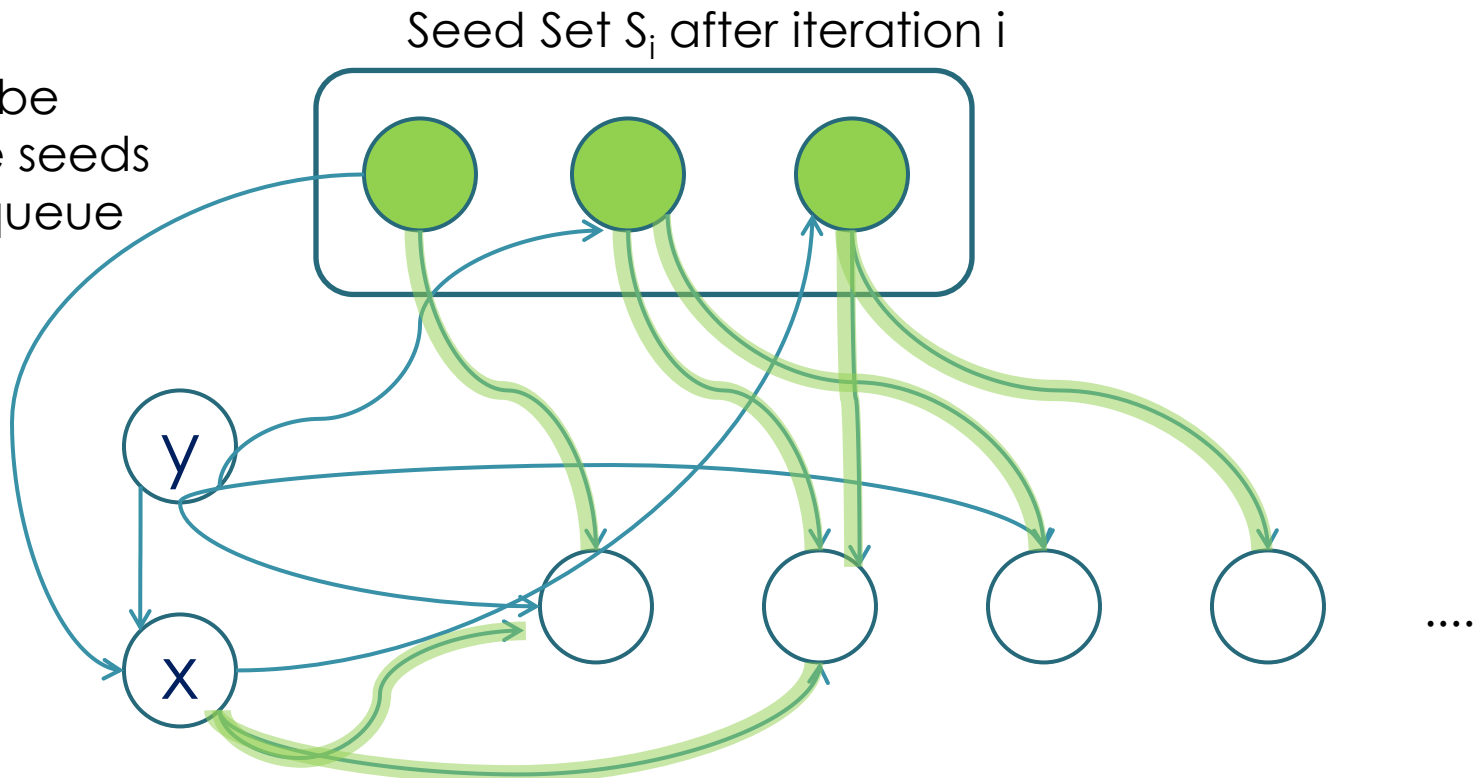
- As the seed set grows, the time spent in estimating spread increases.
 - More paths to enumerate.
- A lot of paths are repeated though.
- The optimization avoids this repetition intelligently.
- A look ahead parameter ' l '.



Look Ahead Optimization (2/2)

$l = 2$ here

Let y and x be prospective seeds from CELF queue



A lot of paths are enumerated repeatedly

$$\begin{aligned} 1. & \sigma(S_i + x) = \text{ad achieved by } S_i + y \\ 2. & \sigma(S_i + x) = \text{ad achieved by } S_i + x \end{aligned}$$

$$\sigma^{V-S_i}(x) + \sigma^{V-x}(S_i)$$

SimPath vs LDAG

TABLE II
SIMPETH'S IMPROVEMENT OVER LDAG

Dataset	Improvement in		
	Spread	Running Time	Memory
NetHEPT	8.7%	21.7%	62.9%
Last.fm	1.7%	42.9%	86.5%
Flixster	8.9%	33.6%	87.5%
DBLP	2.3%	67.2%	87.1%



Other means of speedup

- Community-based approach

[Wang et al. KDD 2010]

- Network sparsification

[Mathioudakis et al, KDD 2011]

- Simulated Annealing

[Jiang et al. AAAI 2011]



Community-based influence maximization

- Use influence parameters to partition graph into communities
 - different from pure structure-based partition
- Influence maximization within each community
- Dynamic programming for selecting top seeds
- Orthogonal with scalable influence computation algorithms

Sparsification of influence networks

- Remove less important edges for influence propagation
- Use action traces to guide graph trimming
- Orthogonal to scalable influence computation algorithms

Simulated annealing

- Following simulated annealing framework
 - not going through all nodes to find a next seed as in greedy
 - find a random replacement (guided by some heuristics)
- Orthogonal to scalable influence computation algorithms

Key takeaways for scalable influence maximization

- tackle from multiple angles
 - reduce #spread evaluations: lazy-forward
 - scale up spread computation:
 - use local influence region (LDAG, MIA, SimPath)
 - use efficient graph structure (trees, DAGs)
 - model specific optimizations (simple path enumerations)
 - graph sparsifications, community partition, etc.
- upshot: can scale up to graphs with millions of nodes and edges, on a **single** machine