

SCALABLE INFLUENCE ESTIMATION IN CONTINUOUS-TIME DIFFUSION NETWORKS

Nan Du ¹

Joint work with Le Song ¹, Manuel Gomez Rodriguez² and Hongyuan Zha¹

¹Georgia Institute of Technology

²Max Planck Institute for Intelligent Systems

MOTIVATION

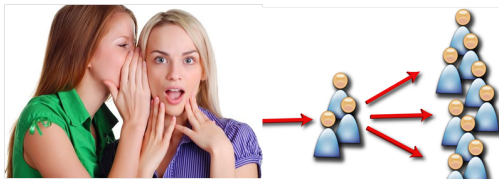
- Networks abstract interactions among entities (media sites, people, organizations, etc.)
- Propagation of news, product reviews, virus, etc. takes place over
 - Information networks
 - Social networks
 - Traffic networks
 - Communication networks

Propagation traces can be extracted from various data sources.

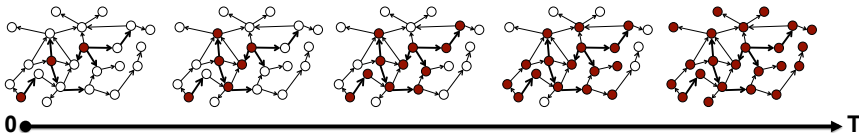


MOTIVATION

- Question : If a new piece of information is released in a few nodes, can it spread, in 1 month, to a million nodes ?
- Question : How can we optimize the selection of the earlier nodes to trigger, *within a time window T* , the largest expected number of follow-ups ?

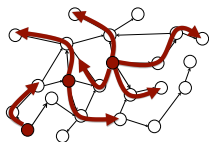


time-sensitive viral marketing



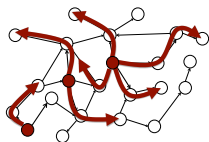
OUTLINE

1 Continuous-time diffusion process.

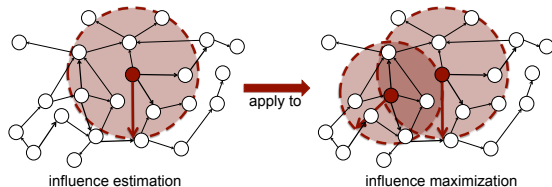


OUTLINE

1 Continuous-time diffusion process.

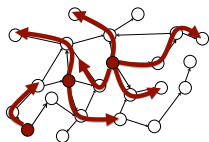


2 Efficient influence estimation and maximization.

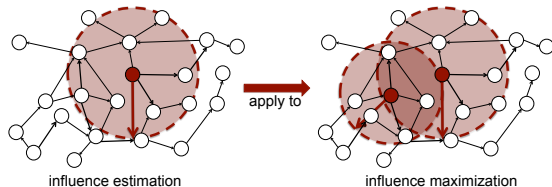


OUTLINE

1 Continuous-time diffusion process.



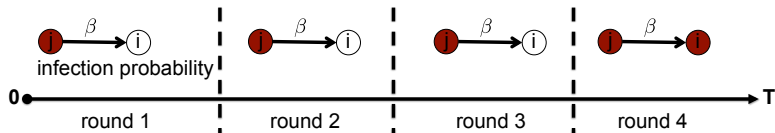
2 Efficient influence estimation and maximization.



3 Experimental evaluation with synthetic and true data.

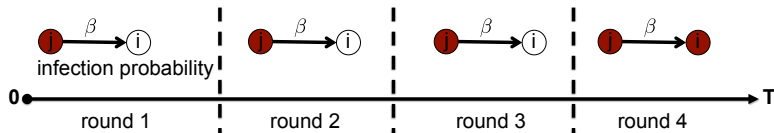
CONTINUOUS VS. DISCRETE TIME DIFFUSION MODEL

- Traditionally, diffusion has been modeled as discrete steps (or rounds).

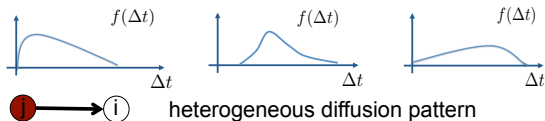


CONTINUOUS VS. DISCRETE TIME DIFFUSION MODEL

- Traditionally, diffusion has been modeled as discrete steps (or rounds).



- However, real time is continuous.

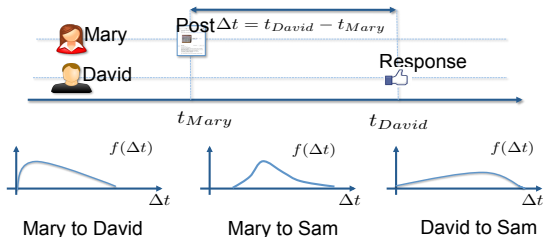


- how long is each round ?
- how to aggregate events within one round ?

CONTINUOUS-TIME INDEPENDENT CASCADE MODEL

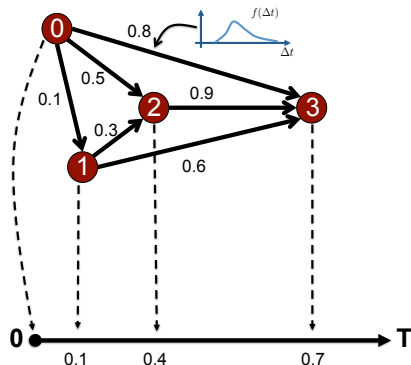
- Node set \mathcal{V} : people, media-sites, organizations, etc.
- Edge set \mathcal{E} : relations, channels, etc.
- Infection : an event occurs to a node, e.g., adopting a product.
- Pairwise conditional density over time

$$f_{ji}(t_j|t_i) = f_{ji}(t_i - t_j)$$



CONTINUOUS-TIME INDEPENDENT CASCADE MODEL

- node 0 is the source $\mathcal{A} = \{0\}$.
- node 0 influences out-going neighbors with $f(\Delta t)$.
- node 1 is infected at $t_1 = 0.1$.
- both node 0 and 1 influence node 2.
- node 1 first infects node 2 since $0.4 < 0.5$.
- node 3 is infected at 0.7 by node 1.



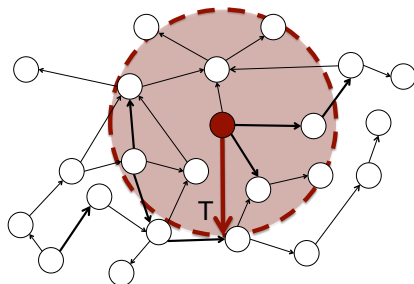
NODE'S PERSPECTIVE

- For each node i , let random variable t_i represent the infection time.

NODE'S PERSPECTIVE

- For each node i , let random variable t_i represent the infection time.
- The influence value of sources \mathcal{A} by time T is

$$\sigma(\mathcal{A}, T) = \mathbb{E} \left[\sum_{i \in \mathcal{V}} \mathbb{I} \{t_i \leq T\} \right] = \sum_{i \in \mathcal{V}} \Pr \{t_i \leq T\}$$



influence estimation

NODE'S PERSPECTIVE

- Directed graphical model representation

$$p(\{t_i\}_{i \in \mathcal{V}}) = \prod_{i \in \mathcal{V}} p(t_i | \{t_j\}_{j \in \pi_i}), \quad \pi_i \text{ is the set of parents}$$

NODE'S PERSPECTIVE

- Directed graphical model representation

$$p(\{t_i\}_{i \in \mathcal{V}}) = \prod_{i \in \mathcal{V}} p(t_i | \{t_j\}_{j \in \pi_i}), \quad \pi_i \text{ is the set of parents}$$

- Marginalization

$$\Pr\{t_i \leq T\} = \int_0^\infty \cdots \int_{t_i=0}^T \cdots \int_0^\infty \left(\prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left(\prod_{j \in \mathcal{V}} dt_j \right)$$

NODE'S PERSPECTIVE

- Directed graphical model representation

$$p(\{t_i\}_{i \in \mathcal{V}}) = \prod_{i \in \mathcal{V}} p(t_i | \{t_j\}_{j \in \pi_i}), \quad \pi_i \text{ is the set of parents}$$

- Marginalization

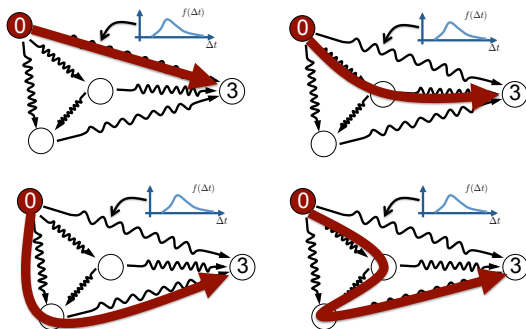
$$\Pr\{t_i \leq T\} = \int_0^\infty \cdots \int_{t_i=0}^T \cdots \int_0^\infty \left(\prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left(\prod_{j \in \mathcal{V}} dt_j \right)$$

- Need to integrate all possible configurations of cascades where $t_i < T$.
- No closed form solution for general heterogeneous transmission function.
- Hard to approximate.

- Mutually independent transmission times $\tau_{ji} = t_i - t_j$

EDGE'S PERSPECTIVE

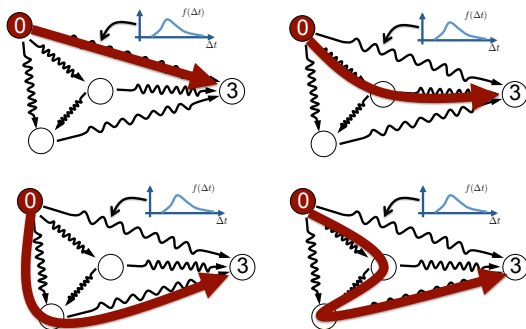
- Mutually independent transmission times $\tau_{ji} = t_i - t_j$
- A network with stochastic edge weights (inter-infection time)



shortest path property

EDGE'S PERSPECTIVE

- Mutually independent transmission times $\tau_{ji} = t_i - t_j$
- A network with stochastic edge weights (inter-infection time)



shortest path property

- t_3 equals to the length of the shortest path from t_0 .

NODE'S VS. EDGE'S PERSPECTIVE

- Node's perspective

$$\sigma(\mathcal{A}, T) = \sum_{i \in \mathcal{V}} \Pr \{t_i \leq T\}$$

$$\Pr \{t_i \leq T\} = \int_0^\infty \cdots \int_{t_i=0}^T \cdots \int_0^\infty \left(\prod_{j \in \mathcal{V}} p(t_j | \{t_l\}_{l \in \pi_j}) \right) \left(\prod_{j \in \mathcal{V}} dt_j \right)$$

- Edge's perspective

$$\sigma(\mathcal{A}, T) = \mathbb{E}_G \left[\sum_{i \in \mathcal{V}} \mathbb{I} \{t_i \leq T\} \right]$$

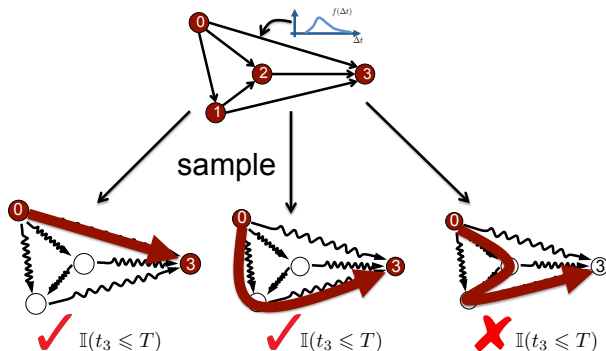
$$G \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

NAIVE SIMULATION

- Given G , t_i is the length of the shortest path.
- Check whether $t_i \leq T$ on many samples.

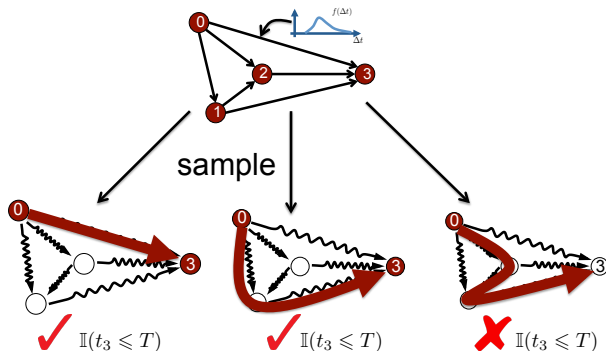
NAIVE SIMULATION

- Given G , t_i is the length of the shortest path.
- Check whether $t_i \leq T$ on many samples.



NAIVE SIMULATION

- Given G , t_i is the length of the shortest path.
- Check whether $t_i \leq T$ on many samples.



- Average the counts across n samples.

$$\sigma(\mathcal{A}, T) \approx \frac{1}{n} \left(\sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_1\} + \dots + \sum_{i \in \mathcal{V}} \mathbb{I}\{t_i \leq T | G_n\} \right)$$

NAIVE SIMULATION

- Using shortest path is not scalable.

NAIVE SIMULATION

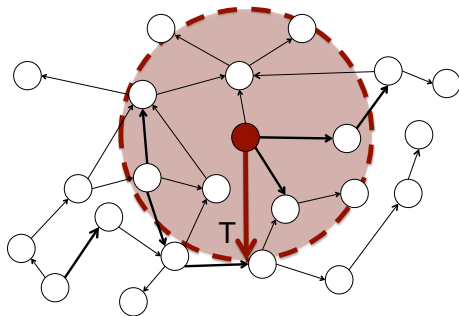
- Using shortest path is not scalable.
- Influence Estimation of a single source j
 - $\sigma(\{j\}, T)$
 - Compute all shortest paths from j to the other nodes.

NAIVE SIMULATION

- Using shortest path is not scalable.
- Influence Estimation of a single source j
 - $\sigma(\{j\}, T)$
 - Compute all shortest paths from j to the other nodes.
- Which source is the best ?
 - Chose j with the largest $\sigma(\{j\}, T)$
 - Try source $j = 0, \dots, |\mathcal{V}| - 1$, $O(|\mathcal{V}|^2)$
- Quadratic in network size
Can not deal with large networks !

NEIGHBORHOOD SIZE ESTIMATION

- Given a sampled network G and source node j , estimate $|\mathcal{N}(\{j\}, T)| = |\{i : t_i \leq T\}|$ the size of neighborhood within distance T .



influence estimation

NEIGHBORHOOD SIZE ESTIMATION

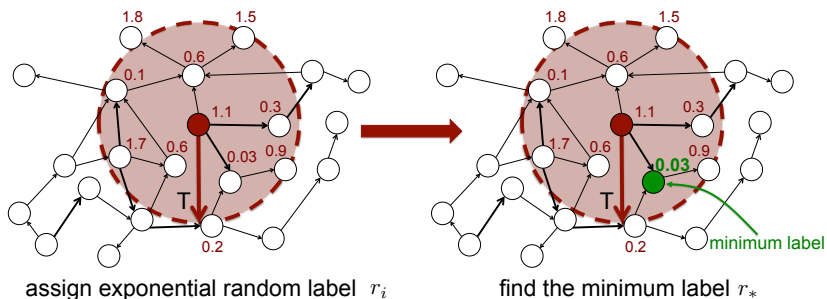
KEY FACT

Given a set of n i.i.d random variable $X^i \sim e^{-x}$, the minimum $X_ \sim ne^{-nx}$.*

NEIGHBORHOOD SIZE ESTIMATION

KEY FACT

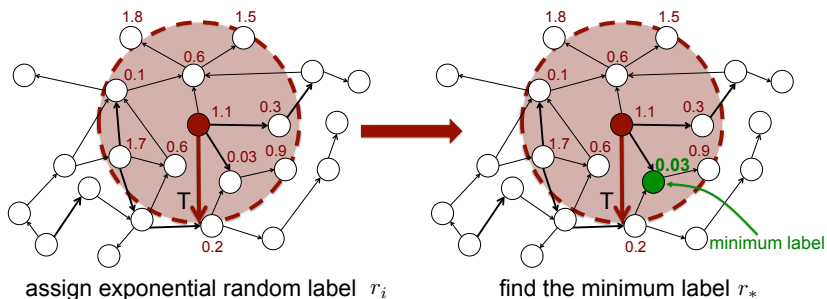
Given a set of n i.i.d random variable $X^i \sim e^{-x}$, the minimum $X_* \sim ne^{-nx}$.



NEIGHBORHOOD SIZE ESTIMATION

KEY FACT

Given a set of n i.i.d random variable $X^i \sim e^{-x}$, the minimum $X_* \sim ne^{-nx}$.

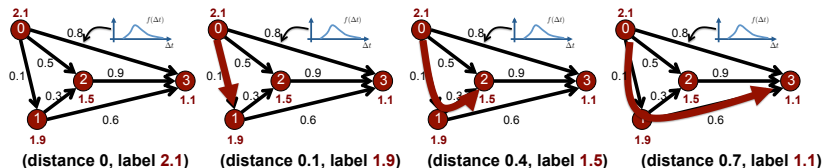


- Find m such least labels, $\{r_*^u\}_{u=1}^m$ to estimate

$|\mathcal{N}(\{j\}, T)| \approx \frac{m-1}{\sum_{u=1}^m r_*^u}$, convert counting problem to estimation problem !

SINGLE SOURCE

- Each node holds a least-label list



Node 0's least-label list

- (Cohen 97) gives the smart algorithm to calculate the least-label list for each node in $\tilde{O}(|\mathcal{E}|)$.
- Estimator $|\mathcal{N}(\{j\}, T)| \approx \frac{m-1}{\sum_{u=1}^m r_u^*}$ is unbiased with variance $O(\frac{1}{m-2})$.
- Sample $m \ll \min(|\mathcal{V}|, |\mathcal{E}|)$ to select the single best source **linearly** in network size !

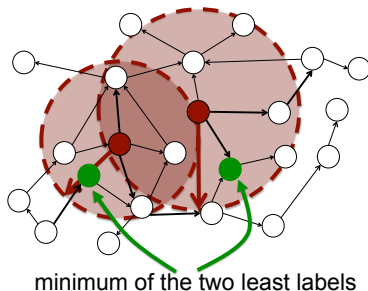
MULTIPLE SOURCES

- Multiple sources \mathcal{A}

$$\mathcal{N}(\mathcal{A}, T) = \bigcup_{s \in \mathcal{A}} \mathcal{N}(s, T).$$

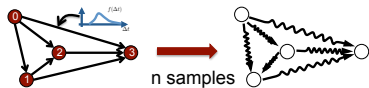
- Reuse least-label list for each single source $s \in \mathcal{A}$

$$r_* = \min_{i \in \mathcal{A}} \min_{j \in \mathcal{N}(i, T)} r_j$$



OVERALL ALGORITHM CONTINEST

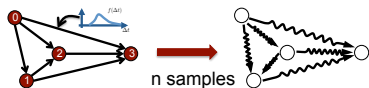
1. Sample n sets of random transmission times



$$\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

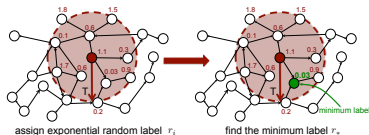
OVERALL ALGORITHM CONTINEST

1. Sample n sets of random transmission times



$$\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

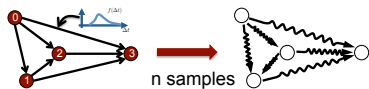
2. Given a set of $\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}}$, sample m sets of random labels



$$\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$$

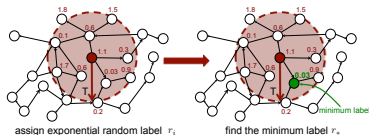
OVERALL ALGORITHM CONTINEST

1. Sample n sets of random transmission times



$$\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}} \sim \prod_{(j,i) \in \mathcal{E}} f_{ji}(\tau_{ji})$$

2. Given a set of $\{\tau_{ij}^l\}_{(j,i) \in \mathcal{E}}$, sample m sets of random labels



$$\{r_i^u\}_{i \in \mathcal{V}} \sim \prod_{i \in \mathcal{V}} \exp(-r_i)$$

3. Estimate $\sigma(\mathcal{A}, T)$ by sample averages

$$\sigma(\mathcal{A}, T) \approx \frac{1}{n} \sum_{l=1}^n \left((m-1) / \sum_{u_l=1}^m r_*^{u_l} \right)$$

THEOREM

Draw the following number of samples for the set of random transmission times

$$n \geq \frac{C\Lambda}{\epsilon^2} \log \left(\frac{2|\mathcal{V}|}{\delta} \right),$$

where Λ depends on \mathcal{A} and T , and for each set of random transmission times, draw m set of random labels. Then $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$ uniformly for all \mathcal{A} with $|\mathcal{A}| \leq C$, with probability at least $1 - \delta$.

THEOREM

Draw the following number of samples for the set of random transmission times

$$n \geq \frac{C\Lambda}{\epsilon^2} \log \left(\frac{2|\mathcal{V}|}{\delta} \right),$$

where Λ depends on \mathcal{A} and T , and for each set of random transmission times, draw m set of random labels. Then $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$ uniformly for all \mathcal{A} with $|\mathcal{A}| \leq C$, with probability at least $1 - \delta$.

- Implications : influence of larger source set \mathcal{A} at the longer time window T requires more samples in the worst case.

THEOREM

Draw the following number of samples for the set of random transmission times

$$n \geq \frac{C\Lambda}{\epsilon^2} \log \left(\frac{2|\mathcal{V}|}{\delta} \right),$$

where Λ depends on \mathcal{A} and T , and for each set of random transmission times, draw m set of random labels. Then $|\hat{\sigma}(\mathcal{A}, T) - \sigma(\mathcal{A}, T)| \leq \epsilon$ uniformly for all \mathcal{A} with $|\mathcal{A}| \leq C$, with probability at least $1 - \delta$.

- Implications : influence of larger source set \mathcal{A} at the longer time window T requires more samples in the worst case.
- In practice : small $m = 5$ achieves good performance. Inaccuracy is canceled out due to large outer-loop n samples.

INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$ is a non-negative, monotonic, submodular function.

INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$ is a non-negative, monotonic, submodular function.
- Greedy algorithm achieves at least a fraction $(1 - 1/e)$ of the optimal value (OPT)

INFLUENCE MAXIMIZATION

- We seek to solve

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq C} \sigma(\mathcal{A}, T)$$

which is NP-hard in general.

- $\sigma(\mathcal{A}, T)$ is a non-negative, monotonic, submodular function.
- Greedy algorithm achieves at least a fraction $(1 - 1/e)$ of the optimal value (OPT)

THEOREM

Suppose the influence $\sigma(\mathcal{A}, T)$ for all \mathcal{A} with $|\mathcal{A}| \leq C$ are estimated uniformly with error ϵ and confidence $1 - \delta$, the greedy algorithm returns a set of sources $\hat{\mathcal{A}}$ such that $\sigma(\hat{\mathcal{A}}, T) \geq (1 - 1/e)\text{OPT} - 2C\epsilon$ with probability at least $1 - \delta$.

EXPERIMENTAL EVALUATION

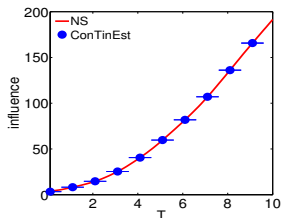
- Synthetic dataset
 - Generate network structure.
 - Weibull pairwise transmission function with randomly chosen parameters.
 - Accuracy of estimated influence (compared to simulation).
 - Quality of selected sources.
 - Scalability.

EXPERIMENTAL EVALUATION

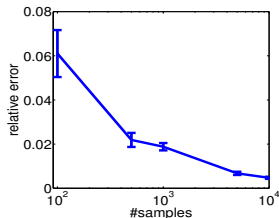
- Synthetic dataset
 - Generate network structure.
 - Weibull pairwise transmission function with randomly chosen parameters.
 - Accuracy of estimated influence (compared to simulation).
 - Quality of selected sources.
 - Scalability.
- Real dataset
 - MemeTracker data (172m news articles 08/2009-09/2009).
 - Infer network structures from hyperlink cascade data.
 - Accuracy of estimated influence (compared to real value).
 - Quality of selected sources on real data.

SYNTHETIC DATASET

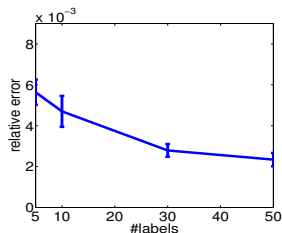
Accuracy of the estimated influence (highest out-degree node)



(a) Influence vs. time



(b) Error vs. #samples

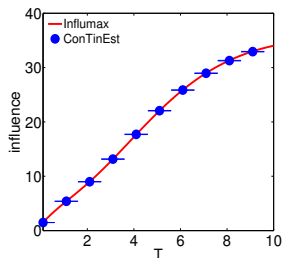


(c) Error vs. #labels

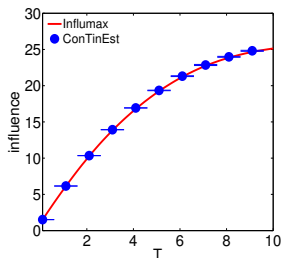
- 100,000 samples for naive simulation (NS).
- $m(\#labels) \ll n(\#samples)$ still achieves good accuracy.
- accuracy does not depend on network structure (1024 nodes, 2048 edges).

SYNTHETIC DATASET

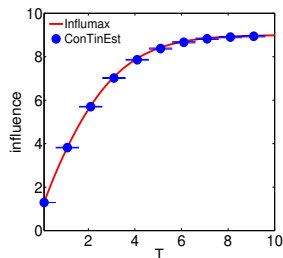
Accuracy of the estimated influence (highest out-degree node)



(a) Core-periphery



(b) Random

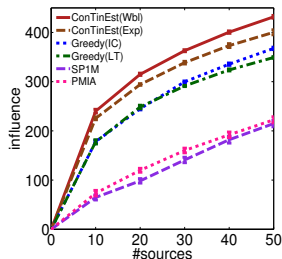


(c) Hierarchical

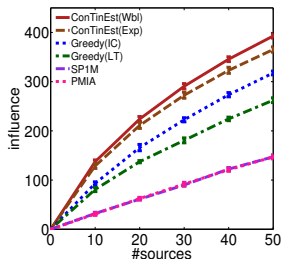
- CONTINEST is close to INFLUMAX (sparse small networks, exponential transmission functions).
- accuracy does not depend on network structure (128 nodes, 141 edges).

SYNTHETIC DATASET

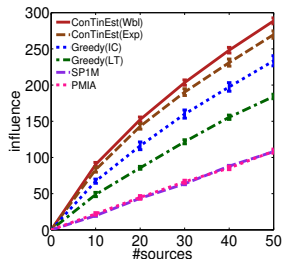
Quality of the selected nodes for influence maximization



(a) Core-periphery



(b) Random

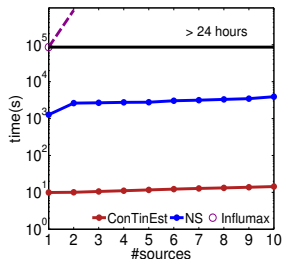


(c) Hierarchical

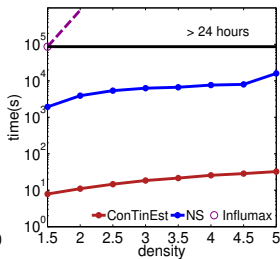
- CONtinEST typically outperforms competitive methods by 20%.
- Performance does not depend on network structure (1024 nodes, 2048 edges).

SYNTHETIC DATASET

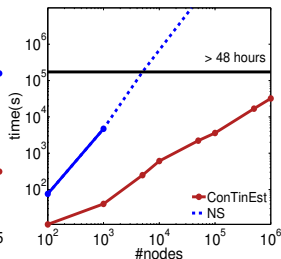
Scalability of influence maximization



(a) # sources
Small network



(b) network density
Small network



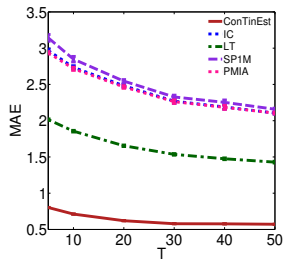
(c) network size
Up to one million nodes

- Small network : 128 nodes.
- Large network : up to 1 million nodes, with density 1.5.
- Our algorithm : sample 10K networks, 5 random labels.

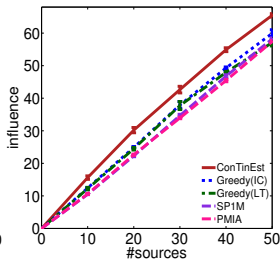
REAL DATASET

- 10,967 hyperlink cascades.
- Use 80% cascades for learning continuous-time diffusion model.
- Select sources based on the learnt model.
- Evaluate influence of the sources using 20% test cascades.
- Compared to discrete-time diffusion models and scalable heuristics.

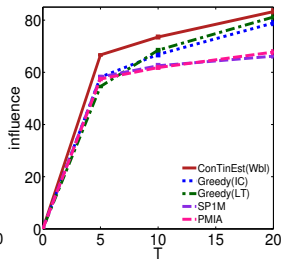
REAL DATASET



(a) Estimation error



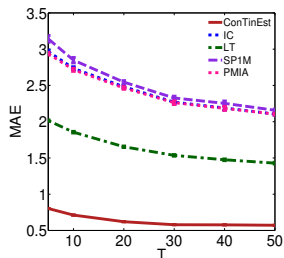
(b) #sources



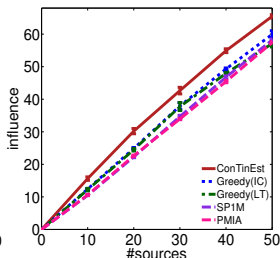
(c) time

- CONtinEST achieves the lowest MAE error.

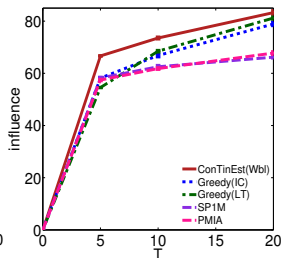
REAL DATASET



(a) Estimation error



(b) #sources



(c) time

- CONTINEST achieves the lowest MAE error.
- CONTINEST produces the set of sources with the largest true influence within short time window.

CONCLUSION

- a novel view of the influence estimation problem in continuous-time diffusion networks.
 - very little assumptions about transmission functions.
 - only depends on temporal cascades induced by diffusion.

CONCLUSION

- a novel view of the influence estimation problem in continuous-time diffusion networks.
 - very little assumptions about transmission functions.
 - only depends on temporal cascades induced by diffusion.
- an efficient randomized algorithm improving :
 - the accuracy of the estimated influence (the lowest MAE in real data).
 - the quality of selected sources (the largest influence within short time period).
 - the scalability (scaling up to millions of nodes in practice).

CONCLUSION

- a novel view of the influence estimation problem in continuous-time diffusion networks.
 - very little assumptions about transmission functions.
 - only depends on temporal cascades induced by diffusion.
- an efficient randomized algorithm improving :
 - the accuracy of the estimated influence (the lowest MAE in real data).
 - the quality of selected sources (the largest influence within short time period).
 - the scalability (scaling up to millions of nodes in practice).
- natural follow-up : product / advertisement allocation with more realistic constraints.