Influence Maximization in Networks

Thanks to Jure Leskovec, Stanford and Panayiotis Tsaparas, Univ. of Ioannina for slides

Agenda

- The Influence Maximization Problem (IMP)
 - (Or, how to create big cascades)
 - (Or, finding the most influential set of nodes)
- IMP Hardness
- IMP Approximation
 - Submodularity
 - Hill Climbing Approximation Algorithm
- IMP Experiments and Remarks

Viral Marketing?

We are more influenced by our friends than strangers



68% of consumers consult friends and family before purchasing home electronics

□50% do research online before purchasing electronics

Viral Marketing

Identify influential customers

Convince them to adopt the product – Offer discount/free samples



Start

These customers endorse the product among their friends

How to Create Big Cascades?

Information epidemics:

- Which are the influential users?
- Which news sites create big cascades?
- Where should we advertise?



VS

Probabilistic Contagion

Independent Cascade Model

- Directed finite G = (V, E)
- Set S starts out with new behavior
 - Say nodes with this behavior are "active"
- Each edge (v, w) has a probability p_{vw}
- If node v is active, it gets <u>one</u> chance to make w active, with probability p_{vw}
 - Each edge fires at most once

Does scheduling matter? No

- *u*, *v* both active, doesn't matter which fires first
- But the time moves in discrete steps

Independent Cascade Model

- Initially some nodes S are active
- Each edge (v, w) has probability (weight) p_{vw}



When node v becomes active:

- It activates each out-neighbor w with prob. p_{vw}
- Activations spread through the network

Most Influential Set of Nodes

- S: is initial active set
- f(S): The expected size of final active set



Set S is more influential if *f(S)* is larger
 f({*a*, *b*}) < *f*({*a*, *c*}) < *f*({*a*, *d*})

Most Influential Set

Problem: (*k* is user-specified parameter)

Most influential set of size k: set S of k nodes producing largest expected cascade size f(S) if activated [Domingos-Richardson '01]



Influence set **X**_a of **a**

Influence set **X_d of d**

• Optimization problem: $\max_{S \text{ of size } k} f(S)$

Why "expected cascade size"? X_a is a result of a random process. So in practice we would want to compute X_a for many realizations and then maximize the "average" value f(S). For now let's ignore this nuisance and simply assume that each node a influences *a* set of nodes X_a



How hard is influence maximization?

Most Influential Subset of Nodes

 Problem: Most influential set of k nodes: set S on k nodes producing largest expected cascade size f(S) if activated
 The optimization problem:

 $\max_{\text{S of size } k} f(S)$

- How hard is this problem?
 - NP-COMPLETE!
 - Show that finding most influential set is at least as hard as a vertex cover

Background: Vertex Cover

Vertex cover problem

(a known NP-complete problem):

Given universe of elements $U = \{u_1, \dots, u_n\}$ and sets $X_1, \dots, X_m \subseteq U$

Are there k sets among X₁,..., X_m such that their union is U?

Goal:

Encode vertex cover as an instance of $\max_{S \text{ of size k}} f(S)$

Influence Maximization is NP-hard

Let a vertex cover instance with sets X₁,..., X_m
 Build a bipartite "X-to-U" graph: Create edge



e.g.: $X_1 = \{u_1, u_2, u_3\}$ Construction: • Create edge $(X_i, u) \forall X_i \forall u \in X_i$ -- directed edge from sets to their elements • Put weight 1 on

• Put weight 1 on each edge (the activation is deterministic)

Vertex Cover as Influence Maximization in X-to-U graph: There exists a set S of size k with f(S)=k+n iff there exists a size k set cover

Note: Optimal solution is always a set of nodes X_i (we never influence nodes "u") This problem is hard in general, could be special cases that are easier.

Summary so Far

Extremely bad news:

- Influence maximization is NP-complete
- Next, good news:
 - There exists an <u>approximation</u> algorithm!
 - For some inputs the algorithm won't find globally optimal solution/set OPT
 - But we will also prove that the algorithm will never do too badly either. More precisely, the algorithm will find a set *S* where *f(S) > 0.63*f(OPT)*, where *OPT* is the globally optimal set.

The Approximation Algorithm

- Consider a <u>Greedy Hill Climbing</u> algorithm to find S:
 - Input:
 - Influence set X_u of each node u: $X_u =$
 - $\{\boldsymbol{v}_1, \boldsymbol{v}_2, \dots\}$
 - If we activate u, nodes $\{v_1, v_2, ...\}$ will eventually get active
 - Algorithm: At each iteration i take the node u that gives best marginal gain: $\max_{u} f(S_{i-1} \cup \{u\})$

 S_i ... Initially active set $f(S_i)$... Size of the union of X_u , $u \in S_i$

(Greedy) Hill Climbing

Algorithm:

- Start with $S_0 = \{ \}$
- For *i* = 1 ... *k*
 - Take node u that max $f(S_{i-1} \cup \{u\})$

• Let
$$\boldsymbol{S_i} = \boldsymbol{S_{i-1}} \cup \{\boldsymbol{u}\}$$

Example:

- Eval. f({a}), ..., f({e}), pick max of them
- Eval. f({d, a}), ..., f({d, e}), pick max
- Eval. f(d, b, a}), ..., f({d, b, e}), pick max



Approximation Guarantee

Claim: Hill climbing produces a solution S
 where: f(S) ≥(1-1/e)*f(OPT) (f(S)>0.63*f(OPT))

[Nemhauser, Fisher, Wolsey '78, Kempe, Kleinberg, Tardos '03]

- Claim holds for functions f(·) with 2 properties:
 - f is monotone: (activating more nodes doesn't hurt) if $S \subseteq T$ then $f(S) \leq f(T)$ and $f({})=0$
 - *f* is submodular: (activating each additional node helps less) adding an element to a set gives less improvement than adding it to one of its subsets: $\forall S \subseteq T$

$$f(S \cup \{u\}) - f(S) \geq f(T \cup \{u\}) - f(T)$$

Gain of adding a node to a small set

Gain of adding a node to a large set

Submodularity– Diminishing returns





Plan: Prove 2 things (1) Our f(S) is submodular (2) Hill Climbing gives nearoptimal solutions (for monotone submodular functions)

Background: Submodular Functions

We must show our *f(·)* is submodular:
∀S ⊆ T

$$f(S \cup \{u\}) - f(S) \geq f(T \cup \{u\}) - f(T)$$

Gain of adding a node to a small set

Gain of adding a node to a large set

Basic fact 1:

• If $f_1(x), ..., f_k(x)$ are submodular, and $c_1, ..., c_k \ge 0$ then $F(x) = \sum_i c_i \cdot f_i(x)$ is also submodular

(Linear combination of submodular functions is a submodular function)

Background: Submodular Functions

$$\forall S \subseteq T: f(S \cup \{u\}) - f(S) \geq f(T \cup \{u\}) - f(T)$$

Gain of adding *u* to a small set Gain of adding *u* to a large set

- Basic fact 2: A simple submodular function
 - Sets X₁, ..., X_m
 - $f(S) = \left| igcup_{k \in S} X_k \right|$ (size of the union of sets X_k , $k \in S$)
 - Claim: f(S) is submodular!



The more sets you already have the less new area a given set will cover



Proof strategy:

- We will argue that influence maximization is an instance of the set cover problem:
 - *f(S)* is the size of the union of nodes influenced by set S



- Note *f(S)* is "random" (a result of a random process) so we need to be careful
- Principle of deferred decision to the rescue!



a

Principle of deferred decision:

- Flip all the coins at the beginning and record which edges fire successfully
- Now we have a deterministic graph!
- Def: Edge is <u>live</u> if it fired successfully
 - That is, we remove edges that did not fire
- What is influence set X_u of node u?
 The set reachable by live-edge paths from u

Influence sets for realization *i*: $X_a^i = \{a, f, c, g\}$ $X_b^i = \{b, c\},$ $X_c^i = \{c\}$ $X_d^i = \{d, e, h\}$



a

- What is an influence set X_u?
 The set reachable by live-edge paths from u
 What is now f(S)?
 - *f_i(S)* = size of the set reachable by live-edge paths from nodes in *S*

For the i-th realization of coin flips

• $f_i(S = \{a, b\}) = |\{a, f, c, g\} \cup \{b, c\}| = 5$ • $f_i(S = \{a, d\}) = |\{a, f, c, g\} \cup \{d, e, h\}| = 7$

Influence sets for realization *i*: $X_a^i = \{a, f, c, g\}$ $X_b^i = \{b, c\},$ $X_c^i = \{c\}$ $X_d^i = \{d, e, h\}$





Activate edges by coin flipping

- Fix outcome $i \in I$ of coin flips
- Xⁱ_v = set of nodes reachable from
 v on live-edge paths
- *f_i(S)* = size of cascades from *S* given coin flips *i*
- $f_i(S) = \left| \bigcup_{v \in S} X_v^i \right| \Rightarrow f_i(S)$ is submodular!

• X_{v}^{i} are sets, $f_{i}(S)$ is the size of their union

- Expected influence set size:
 - $f(S) = \sum_{i \in I} f_i(S) \Rightarrow f(S)$ is submodular!
 - f(S) is a linear combination of submodular functions





RECAP: Influence Maximization

Find most influential set S of size k: largest expected cascade size *f(S)* if set *S* is activated



Network, each edge activates with prob. **p**₁₁

Want to solve:

$$\max_{|S|=k} f(S) = \sum_{i \in I} f_i(S)$$

Consider **S={a,d}** then: $f_1(S)=5, f_2(S)=4, f_3(S)=3$ and **f(S) = 12**

Activate edges by coin flipping Multiple realizations *i*:



... influence set of node d

Plan: Prove 2 things (1) Our f(S) is submodular (2) Hill Climbing gives nearoptimal solutions (for monotone submodular functions)

Proof for Hill Climbing

<u>Claim:</u> If f(S) is monotone and submodular. Hill climbing produces a solution S where: $f(S) \ge (1 - \frac{1}{e}) \cdot f(OPT)$ In other words: $f(S) \ge 0.63 \cdot f(OPT)$

The setting:

- Keep adding nodes that give the largest gain
- Start with $S_0 = \{\}$, produce sets S_1, S_2, \dots, S_k
- Add elements one by one
- Let $OPT = \{t_1 \dots t_k\}$ be the optimal set (OPT) of size k
- We need to show: $f(S) \ge (1 \frac{1}{e}) f(OPT)$

Proof Overview

Define: Marginal gain: δ_i = f(S_i) - f(S_{i-1})
Proof: 3 steps:

• 0) Lemma: $f(A \cup B) - f(A) \le \sum_{j=1}^{k} [f(A \cup \{b_j\}) - f(A)]$ • where: $B = \{b_1, \dots, b_k\}$ and $f(\cdot)$ is submodular • 1) $\delta_{i+1} \ge \frac{1}{k} [f(OPT) - f(S_i)]$ • 2) $f(S_{i+1}) = (1 - \frac{1}{k}) f(S_i) + \frac{1}{k} f(OPT)$ • 3) $f(S_k) \ge (1 - \frac{1}{e}) f(OPT)$

Step zero: Basic Hill Climbing Fact

- $f(A \cup B) f(A) \le \sum_{j=1}^{k} [f(A \cup \{b_j\}) f(A)]$
 - where: $B = \{b_1, \dots, b_k\}$ and $f(\cdot)$ is submodular
- Proof:
- Let $B_i = \{b_1, ..., b_i\}$, so we have $B_1, B_2, ..., B_k (= B)$ • $f(A \cup B) - f(A) = \sum_{i=1}^{k} [f(A \cup B_i) - f(A \cup B_{i-1})]$ • = $\sum_{i=1}^{k} [f(A \cup B_{i-1} \cup \{b_i\}) - f(A \cup B_{i-1})]$ $\leq \sum_{i=1}^{k} [f(A \cup \{b_i\}) - f(A)]$ Work out the sum. Everything but 1st and last term cancel out: $f(A \cup B_1) - f(A \cup B_0)$ $+ f(A \cup B_2) - f(A \cup B_1)$ By submodularity since $A \cup X \cup \{b\} \supseteq A \cup \{b\}$ $+f(A \cup B_3) - f(A \cup B_2) \dots$ $+ f(A \cup B_k) - f(A \cup B_{k-1})$

Step one: What is δ_i gain at step *i*?

Remember: $\delta_i = f(S_i) - f(S_{i-1})$

• $f(OPT) \le f(S_i \cup OPT)$ (by monotonicity) $= f(S_i \cup OPT) - f(S_i) + f(S_i)$ $\leq \sum_{j=1}^{k} \left[f\left(S_i \cup \{t_j\}\right) - f(S_i) \right] + f(S_i)$ (by prev. slide) $\leq \sum_{i=1}^{k} [\delta_{i+1}] + f(S_i)$ $OPT = \{ t_1, \dots, t_k \}$ t_i is j-th element of the optimal solution. $= f(S_i) + k \,\delta_{i+1}$ Rather than choosing t_i let's greedily choose the **best element** q_i, which • Thus: $f(OPT) \leq f(S_i) + k \delta_{i+1}$ gives a gain of δ_{i+1} . So, $f(S_i \cup \{t_j\}) \leq \delta_{i+1}$. This is the "hill-climbing" $\Rightarrow \delta_{i+1} \geq \frac{1}{\nu} [f(OPT) - f(S_i)]$ assumption.

Step two: What is f(S_{i+1})?

- We just showed: $\delta_{i+1} \ge \frac{1}{k} [f(OPT) f(S_i)]$
- What is *f*(*S*_{*i*+1})?

•
$$f(S_{i+1}) = f(S_i) + \delta_{i+1}$$

 $\bullet \ge f(S_i) + \frac{1}{k} [f(OPT) - f(S_i)]$

$$\bullet = \left(1 - \frac{1}{k}\right)f(S_i) + \frac{1}{k}f(OPT)$$

• What is $f(S_k)$?

Step three: What is f(S_k)?

• Claim:
$$f(S_i) \ge \left[1 - \left(1 - \frac{1}{k}\right)^i\right] f(OPT)$$

Proof by induction:

• *i* = 0:

•
$$f(S_0) = f(\{\}) = 0$$

• $\left[1 - \left(1 - \frac{1}{k}\right)^0\right] f(OPT) = 0$

Step three: What is $f(S_k)$?

- Given that this is true for S_i : $f(S_i) \ge \left| 1 \left(1 \frac{1}{k}\right)^i \right| f(OPT)$
- **Proof by induction:** • At *i* + 1:

•
$$f(S_{i+1}) \ge \left(1 - \frac{1}{k}\right) f(S_i) + \frac{1}{k} f(OPT)$$

• $\ge \left(1 - \frac{1}{k}\right) \left[1 - \left(1 - \frac{1}{k}\right)^i\right] f(OPT) + \frac{1}{k} f(OPT)$
the claim
• $= \left[1 - \left(1 - \frac{1}{k}\right)^{i+1}\right] f(OPT)$
Two slides ago we showed:
 $f(S_{i+1}) \ge \left(1 - \frac{1}{k}\right) f(S_i) + \frac{1}{k} f(OPT)$

What is f(S_k)?

• Thus: $f(S) = f(S_k) \ge \left[1 - \left(1 - \frac{1}{k}\right)^k\right] f(OPT)$ $\leq \frac{1}{e}$ • So:

$$f(S_k) \ge \left(1 - \frac{1}{e}\right) f(OPT)$$

qed.

Solution Quality

We just proved:

Hill climbing finds solution S which
 f(S) ≥ (1-1/e)*f(OPT) i.e., f(S) ≥ 0.63*f(OPT)

This is a data independent bound

- This is a worst case bound
- No matter what is the input data, we know that the Hill-Climbing will never do worse than 0.63*f(OPT)

Evaluating f(S)?

How to evaluate f(S)?

Still an open question of how to compute it efficiently

But: Very good estimates by simulation

- Repeating the diffusion process often enough (polynomial in *n*; 1/ε)
- Achieve (1±ε)-approximation to f(S)
- Generalization of Nemhauser-Wolsey proof: Greedy algorithm is now a (1-1/e- ε')approximation

Experiments and Concluding Thoughts

Experiment Data

- A collaboration network: co-authorships in papers of the arXiv high-energy physics theory:
 - 10,748 nodes, 53,000 edges
 - Example cascade process: Spread of new scientific terminology/method or new research area
- Independent Cascade Model:
 - Case 1: Uniform probability p on each edge
 - Case 2: Edge from v to w has probability
 1/deg(w) of activating w.

Experiment Settings

- Simulate the process 10,000 times for each targeted set
 - Every time re-choosing edge outcomes randomly

Compare with other 3 common heuristics

- Degree centrality: Pick nodes with highest degree
- Distance centrality: Pick nodes in the "center" of the network
- Random nodes: Pick a random set of nodes

Independent Cascade Model



Uniform edge firing probability *p*_{uv}

Independent Cascade Model



• **Notice:** Greedy approach is slow!

- For a given network G, repeat 10,000s of times:
 - Flip coin for each edge and determine influence sets under coin-flip realization i
 - Each node u is associated with 10,000s influence sets X_uⁱ
- Greedy's complexity is $O(k \cdot n \cdot R \cdot M)$
 - n ... number of nodes in G
 - k ... number of nodes be selected/influenced
 - R ... number of simulation rounds
 - *m* ... number of edges in *G*

Cottage Industry of Heuristics

 Many researchers have since proposed heuristics that work well in practice and run faster than the greedy algorithm



Figure 2: Influence spreads of different algorithms on the collaboration graph NetPHY under the independent cascade model (n = 37, 154, m = 231, 584, and p = 0.01).



Figure 4: Running times of different algorithms on the collaboration graph NetPHY under the independent cascade model (n = 37, 154, m = 231, 584, p = 0.01, and k = 50).

Open Questions

More realistic marketing:

- Different marketing actions increase likelihood of initial activation, for several nodes at once
- Study more general influence models
 - Find trade-offs between generality and feasibility
- Deal with negative influences
 - Model competing ideas
- Obtain more data (better models) about how activations occur in real social networks