

CS630 Graduate Algorithms

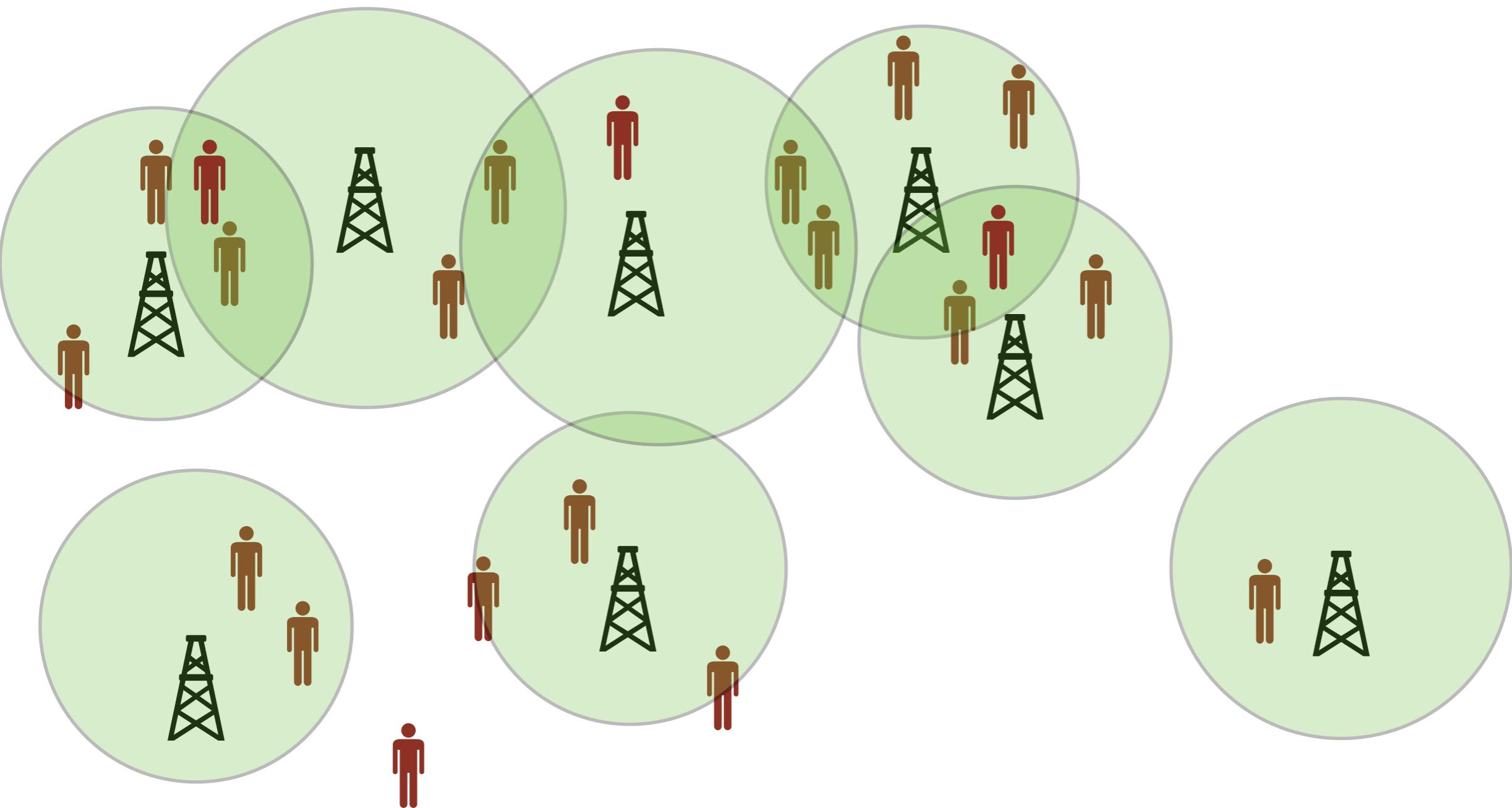
October 1, 2024

by Dora Erdos and Jeffrey Considine

- Monotone submodular functions
 - coverage problems
 - greedy optimization

Radio stations

Each station has a broadcast range. Where to broadcast from to reach the maximum audience, if we only have a budget to broadcast from k stations?



Max k-Coverage Problem

Set Cover: Given a universe $U = \{u_1, u_2, \dots, u_n\}$ of elements and a collection $S = \{S_1, S_2, \dots, S_m\}$ of subsets of U , find a minimum number of the sets in S such that their union contains every item in U .

\hookrightarrow target is fixed

\hookrightarrow not fixed

Max k-Coverage Problem: Given a universe $U = \{u_1, u_2, \dots, u_n\}$ of elements and a collection $S = \{S_1, S_2, \dots, S_m\}$ of subsets of U and an integer k , find k sets in S such that the number of elements covered by their union is maximized.

\curvearrowleft fixed # of sets

\curvearrowright not fixed

Max k-coverage greedy algorithm

Algorithm:

Algorithm 1: GreedySC(U, S_1, \dots, S_m) *additional input: K*

```
1  $X \leftarrow U$  /* uncovered elements in U */  
2  $C \leftarrow$  empty set of subsets;  
3 while  $X$  is not empty do → for  $j = 1 \dots K$ :  
4   Select  $S_i$  that covers the most items in  $X$ ;  
5    $C \leftarrow C \cup S_i$ ;  
6    $X \leftarrow X \setminus S_i$ ;  
7 return  $C$ ;
```

update algo to solve K -coverage

worst # of iters for GreedySC $n \rightarrow$ each iter
only covers
one additional
item

Max k-coverage greedy algorithm

Algorithm: for k iterations select the set that covers the most additional elements.

Algorithm 1: GreedySC(U, S_1, \dots, S_m, k)

```
1  $X \leftarrow U$  /* uncovered elements in U */  
2  $C \leftarrow$  empty set of subsets;  
3 while  $X$  is not empty do For j=1...k do  
4   Select  $S_i$  that covers the most items in  $X$ ;  
5    $C \leftarrow C \cup S_i$ ;  
6    $X \leftarrow X \setminus S_i$ ;  
7 return  $C$ ;
```

Tophat

Set Cover: find the minimum number of sets to cover the universe

Max k-Coverage: find k sets to cover as much of the universe as possible

Select all correct statements:

- A. the greedy algorithm for Max k-Coverage has constant number of iterations
 k is part of the input. *possible inputs* $k = \log n$
- B. for $k = n$ the output for the greedy algorithm for Set Cover and Max k-Coverage
are identical
↓
stop SetCover once we
covered everything $\rightarrow \min$
of sets
- C. the greedy algorithm for Set Cover has a worse approximation ratio than the
algorithm for Max k-Coverage
- D. the greedy algorithm for Max k-Coverage has a constant approximation ratio
while the algorithm for Set Cover depends on n
↑
really cool!!

Max k-coverage approximation

$$\text{calculus: } \left(1 - \frac{1}{t}\right)^t < \frac{1}{e}$$

Theorem: The greedy algorithm has approximation factor $1 - \left(1 - \frac{1}{k}\right)^k > 1 - \frac{1}{e} \approx 63\%$

meaning: suppose: optimal solutions covers $\underline{z^*}$ items
objective: cover as many items as possible

\underline{z} output of greedy:

$$\underline{z^*} \geq \underline{z} \geq 0.63 \underline{z^*}$$

$$\boxed{\begin{array}{l} \text{def: } c \geq 1 \\ \underline{z} \leq \underline{z^*} \leq c \cdot \underline{z} - \\ \underline{z^*} \leq c \cdot \underline{z} \end{array}} \quad c = e \approx 2.71... \quad \frac{\underline{z^*}}{e} \leq \underline{z}$$

remember: for the Set Cover problem, if the optimal solution uses L sets, then the approximation factor of the greedy algorithm is $\underline{L} \cdot \ln n$

What's the difference?

Max k-coverage approximation

Theorem: The greedy algorithm has approximation factor $1 - \left(1 - \frac{1}{k}\right)^k > 1 - \frac{1}{e} \approx 63\%$ meaning:

- the greedy algorithm covers at least ~63% of the items that an optimal cover with k sets would

remember: for the Set Cover problem, if the optimal solution uses L sets, then the approximation factor of the greedy algorithm is $L \cdot \ln n$

What's the difference?

- for the k -coverage problem the approximation has *constant* ratio, for set cover it depends on the *input size* n
- (note that k is not constant, it's part of the input!)
- intuitively the first k sets cover larger ratio of the points, as we select more sets the marginal gain of extra elements covered is diminishing

Max k-coverage approximation

$$\text{calculus: } t > 0 : \left(1 - \frac{1}{t}\right)^t < \frac{1}{e}$$

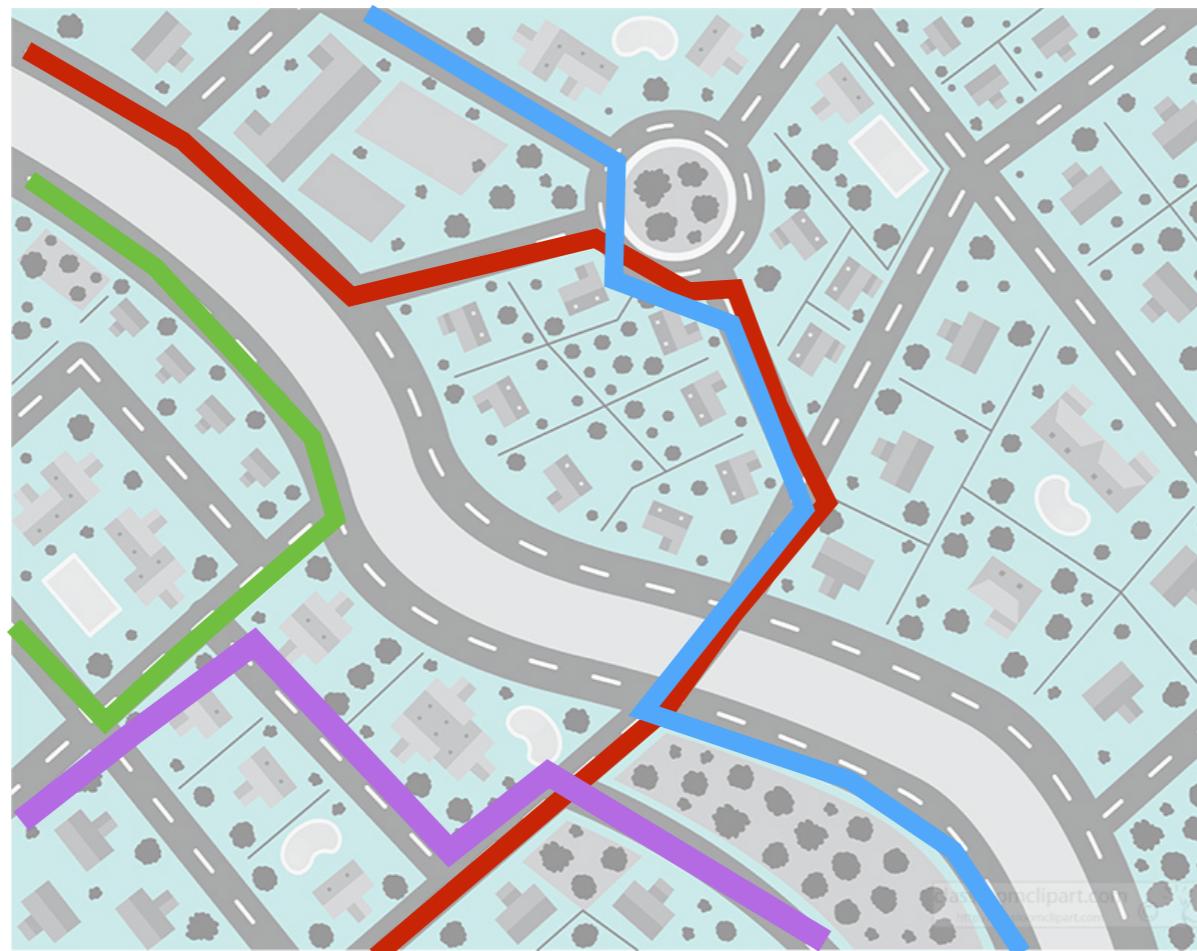
Theorem: the greedy algorithm has approximation factor $1 - \left(1 - \frac{1}{k}\right)^k > 1 - \frac{1}{e} \approx 63\%$

proof:

- Z is the set of items covered by the optimal solution
- among the k sets in the optimal solution, there is at least one set that covers $1/k$ fraction of Z
- since Greedy-k-SC selects the largest set, it also covers at least Z/k items
- after the first iteration at most $Z \left(1 - \frac{1}{k}\right)$ remain uncovered
- since greedy selects the largest marginal gain, it covers at least $1/k$ of the remaining elements in Z in each iteration: $Z \left(1 - \frac{1}{k}\right)^k$
- after k rounds there are at least $Z - Z \left(1 - \frac{1}{k}\right)^k = Z \left(1 - \left(1 - \frac{1}{k}\right)^k\right)$ points covered

Detecting Potholes

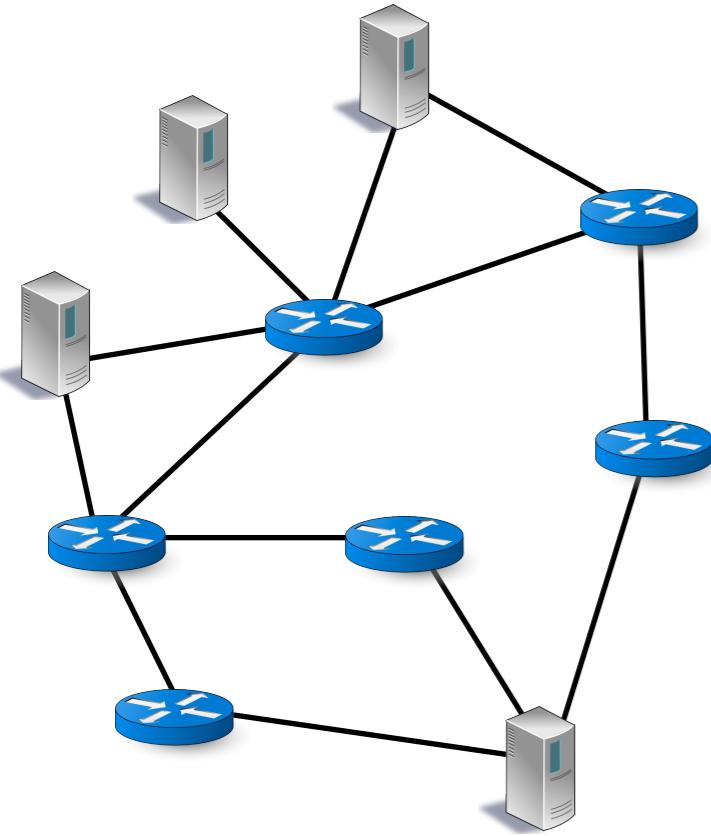
- Mount sensors on busses to detect potholes in the road along their routes
- Bus routes overlap so different routes may cover the same streets
- Given a small budget of sensors, which bus routes should we equip with sensors to detect as many potholes as possible?



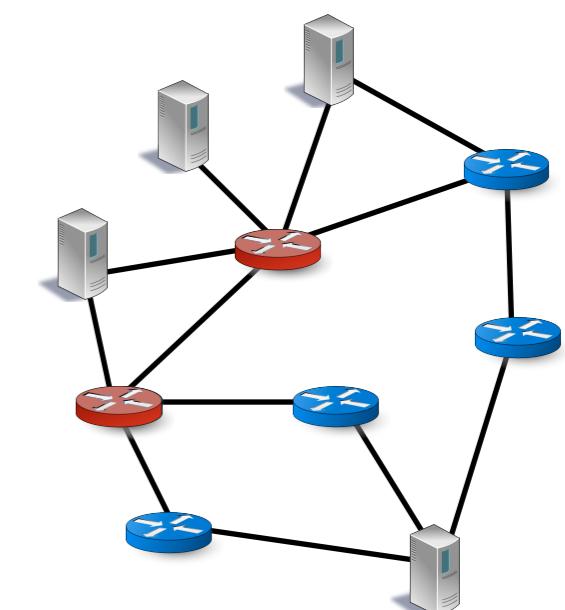
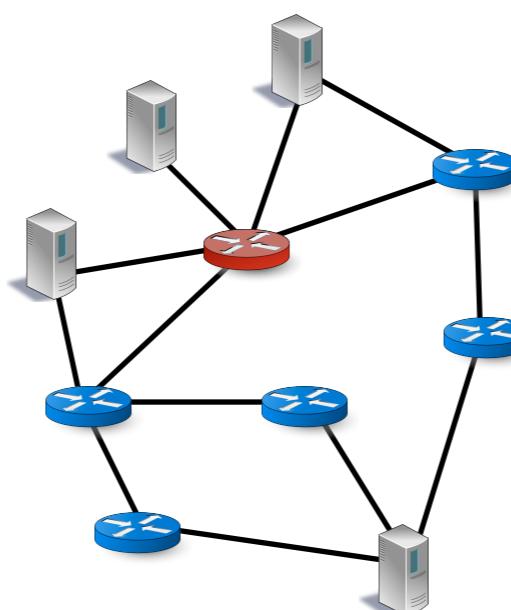
Ali, Dyo [2017]: <https://www.scitepress.org/Papers/2017/64698/64698.pdf>

Covering shortest paths

Select two routers that together cover the most shortest paths.



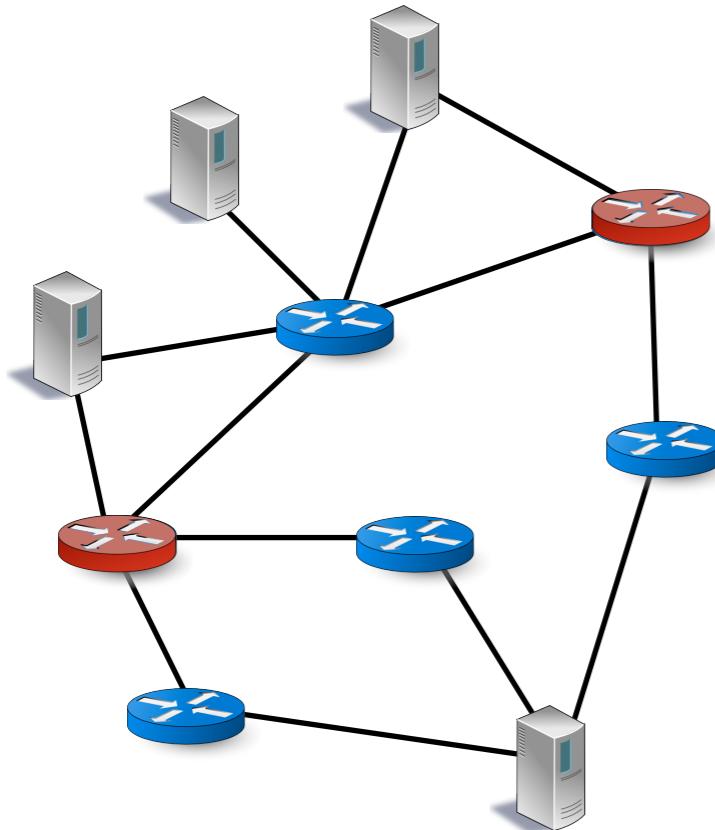
Greedy: select router that covers most additional paths



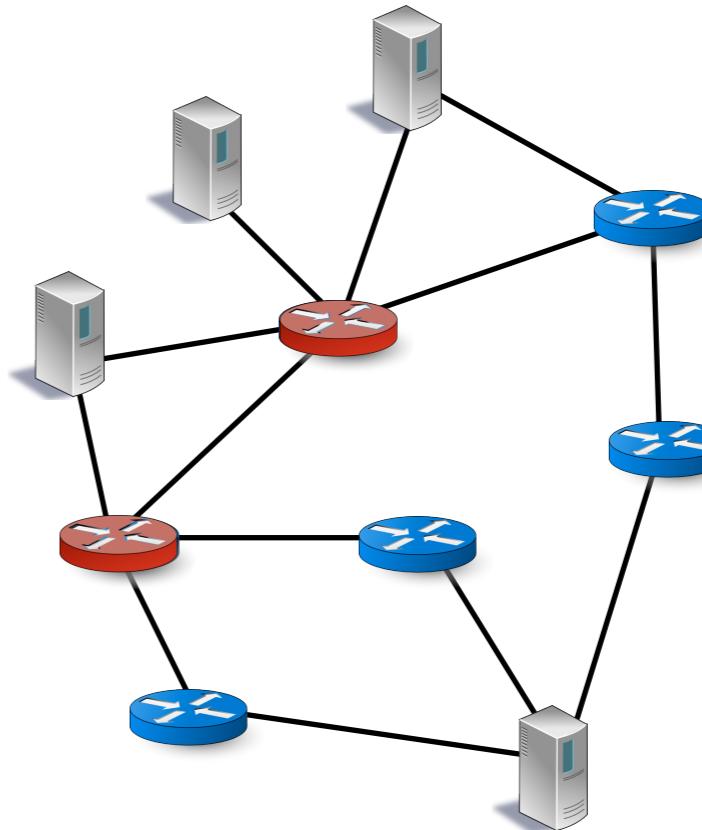
Covering shortest paths

Select **two** routers that together cover the most shortest paths.

Optimal solution



Greedy: select router that covers most additional paths



Bulk Pricing



BULK PRICE
ELIGIBLE **\$3⁸⁵** /piece
Buy 50 or more \$3.27

Model# 769887219614
2 in. x 4 in. x 96 in. Prime Kiln-Dried Whitewood Stud

 Pickup
2,500 in stock at Watertown

 Delivery
Unavailable



BULK PRICE
ELIGIBLE **\$3⁸⁵**
Buy 50 or more \$3.27

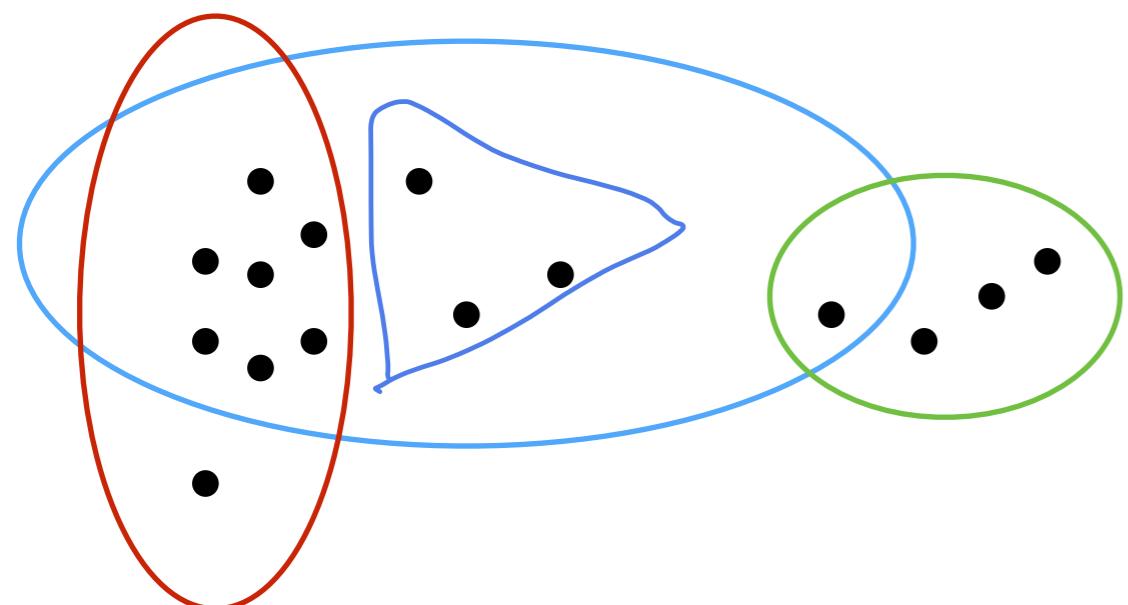
Model# 058449
2 in. x 4 in. x 8 ft. Prime Stud

 Pickup
2,500 in stock at Watertown

 Delivery
Scheduled Delivery

Greedy algorithm for coverage problems

- in each problem we assign some positive value to a set of objects
- diminishing returns
 - the additional benefit of one more set is less as more sets are selected



1. select blue set
2. select green, not red



benefit of adding the
red set to Δ : 8 additional
covered
adding red to Δ : 1 additional
cost.

- natural greedy algorithm:
 - For k iterations repeatedly select the object with largest gain towards our objective function.

Objective function for coverage problems

set function: a function $f: 2^X \rightarrow \mathbb{R}_+$ that takes sets as input and outputs numbers

- 2^X is the set of all subsets of X , think of the set represented as a bit vector

$$X : \begin{array}{ccc} \checkmark & \checkmark & \checkmark \\ \underline{x} & \underline{x} & \underline{x} \\ & 1 & 2 & 3 \end{array} \quad \dots \quad \overset{|X|}{\overbrace{\dots}} \quad |X|$$

build a subset: for each item decide whether part of the set or not

$$2 \cdot 2 \cdot 2 \cdot \dots \cdot 2 = 2^{|X|} \text{ different subset}$$

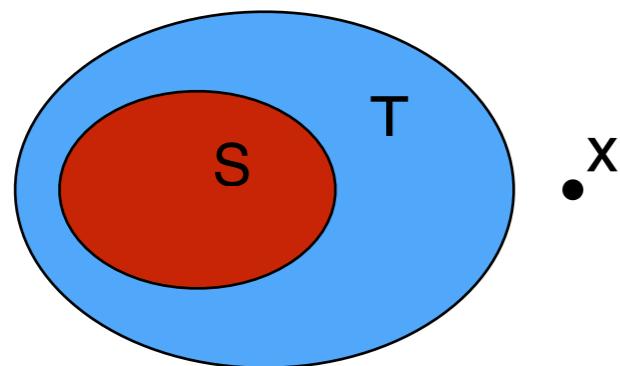
Submodular functions

intuition:

$f(S)$ = how items does S cover

The set function $f: 2^X \rightarrow \mathbb{R}_+$ is **submodular** if for every $S \subset T \subset X$ and $x \in X \setminus T$

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



increase in the value of $f(\cdot)$
is more if the same item
is added to a smaller
subset vs. larger

f is monotone increasing if for every $S \subseteq T$ we have $f(S) \leq f(T)$

marginal benefit of adding smg
to a smaller set is higher

Set Cover:

suppose
• items u_1, u_2, \dots, u_K are currently covered
• items $u_1, u_2, \dots, u_K, u_{K+1}, u_{K+2}, \dots, u_m$

add the same subset S_i to both.

How many additional items are covered?

TopHat - monotone submodular functions

Example:
 $S \subseteq T \subseteq X$

Which of the following functions are monotone submodular? ex1.



or



A. X is the universe, x in X is one specific item, S is a subset of ~~X~~

$$\checkmark f(S) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{otherwise} \end{cases}$$

~~ex1:~~ $f(S \cup \{x\}) - f(S) = 1$
~~ex2:~~ $f(T \cup \{x\}) - f(T) = 1$

~~ex1:~~ $f(S \cup \{x\}) - f(S) = 1$
~~ex2:~~ $f(T \cup \{x\}) - f(T) = 0$

~~B. $f(S)$ = assign a random value to each set as input to the problem~~

no control over what is larger or smaller

~~C. $f(S) = \text{constant}$~~

~~D. $f(S) = \begin{cases} 1 & |S| \text{ is even} \\ 2|S| & S \text{ is odd} \end{cases}$~~

f fluctuate between 1 and 2 x set size with each additional item

~~E. $f(S) = 5|S|$~~ any set adding 1 more item:

$$f(S \cup \{x\}) - f(S) = (|S| + 1) \cdot 5 - |S| \cdot 5 = 5$$

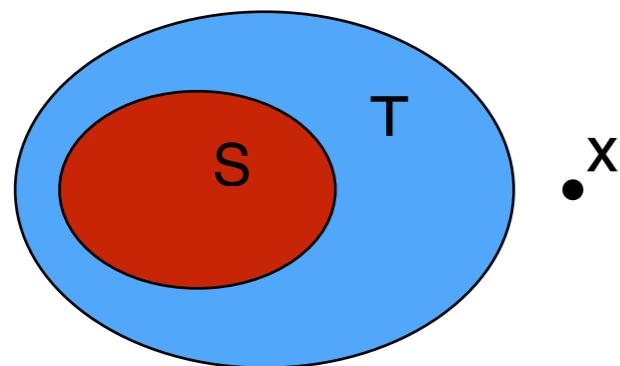
Submodular functions

The set function $f: 2^X \rightarrow \mathbb{R}_+$ is **submodular** if for every $S \subset T \subset X$ and $x \in X \setminus T$

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

part A from TopHat slide:

if x is not part of S
 \Rightarrow adding x increasing
the value
 $f(S \cup \{x\}) = f(S) + 1$



f is **monotone** increasing if for every $S \subseteq T$ we have $f(S) \leq f(T)$

examples of monotone submodular functions:

$$f(S) = c \cdot |S| \quad c > 0$$

$$f(S) = \sum_{i \text{ in } S} w_i \text{ where } w_i \geq 0 \text{ linear functions}$$

$$\text{budget-additive } f(S) = \min\{B, \sum_{i \text{ in } S} w_i\}$$

coverage functions - items, paths, sets, ...

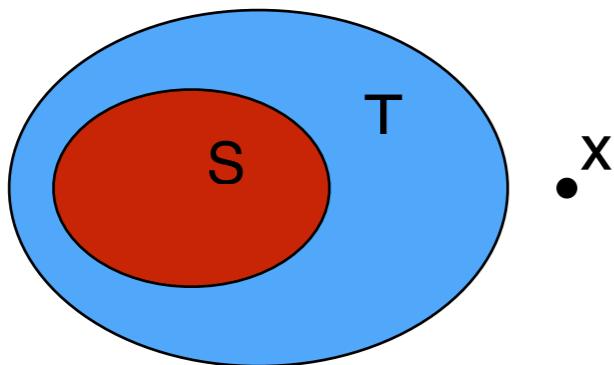
entropy of random variables, information gain

if x is already
in S
 \Rightarrow "adding" x
doesn't change
the set
 \Rightarrow no increase.

Submodular functions

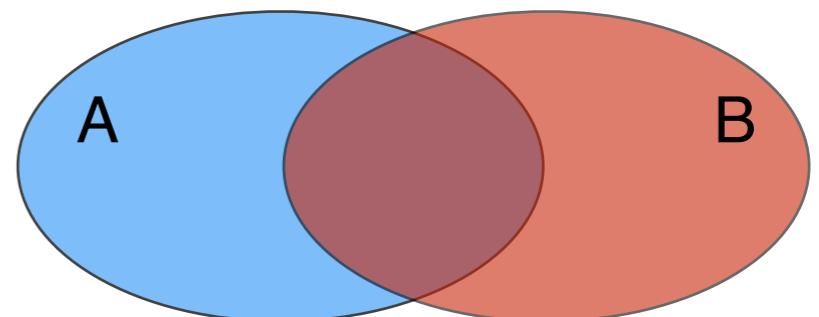
The set function $f: 2^X \rightarrow \mathbb{R}_+$ is submodular if for every $S \subset T \subset X$ and $x \in X \setminus T$

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



Equivalent definition: f is submodular if for every $A, B \subset X$

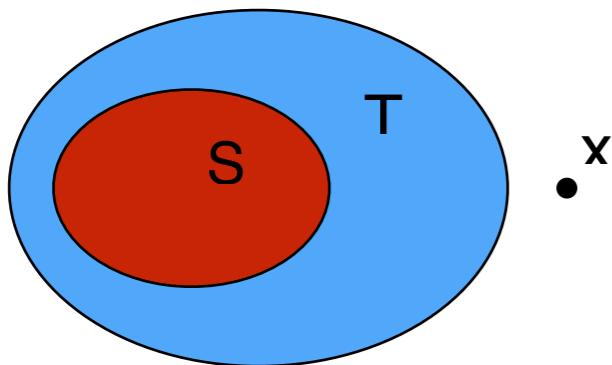
$$f(A \cup B) + f(A \cap B) \leq f(A) + f(B)$$



Submodular functions

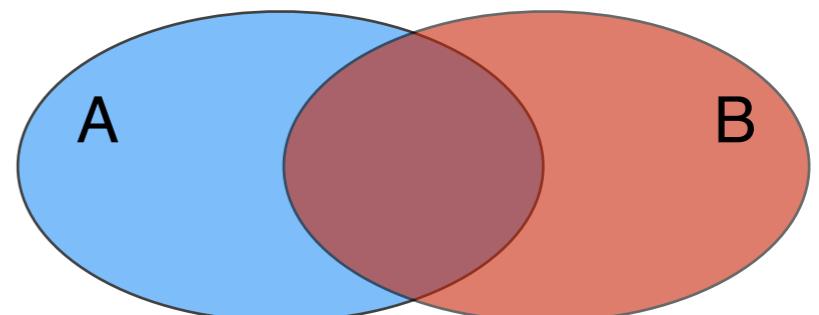
The set function $f: 2^X \rightarrow \mathbb{R}_+$ is submodular if for every $S \subset T \subset X$ and $x \in X \setminus T$

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



Equivalent definition: f is submodular if for every $A, B \subset X$

$$f(A \cup B) + f(A \cap B) \leq f(A) + f(B)$$



proof:

\Rightarrow setting $A = T, B = S \cup \{x\}$ we get the formula

\Leftarrow use $A \cap B \subseteq B$ inductively apply the inequality to each element in $A \setminus B$ to get
 $F(A \cup B) - F(B) \leq F(A) - F(A \cap B)$

Greedy algorithm for monotone submodular functions

Suppose that the objective function f of some maximization problem is *monotone submodular*.

Algorithm 1: GreedySubmodular($X, S_1, \dots, S_m, k f()$)

```
/* X is the universe of elements,  $S_i$  are subsets,  $k$  is an
   int,  $f( )$  is a submodular function */  
1  $C \leftarrow \emptyset$  /*  $C \subseteq X$  currently covered items in  $X$  */  
2 for  $i = 1$  to  $k$  do  
3   | find  $i$  to maximize  $f(C \cup S_i) - f(C)$ ;  
4   |  $C \leftarrow C \cup \{S_i\}$ ;  
5 return  $C$ 
```

submodular functions and complexity

problem type	maximization	minimization
unconstrained	NP-hard some approximations	polynomial via convex optimization
constrained - select k	NP-hard constant approx ratio $\frac{1}{e}$	usually NP-hard to approximate

Greedy approximation factor

Theorem [Nemhauser, Wolsey, Fisher, 1978]: For any maximization problem with a monotone submodular objective function the greedy algorithm yields a $(1 - \frac{1}{e})$ -approximation.

$$1 - \frac{1}{e} \approx 63\%$$

Why is this useful?

Greedy approximation factor

Theorem [Nemhauser, Wolsey, Fisher, 1978]: For any maximization problem with a monotone submodular objective function the greedy algorithm yields a $(1 - 1/e)$ -approximation.

Why is this useful?

- for optimization problems - which are often NP-C - the most simple greedy algorithm is a pretty good optimization.
 - “pretty good” = constant!
- it’s enough to prove that the function the problem is maximizing is indeed monotone submodular.

Product adoption via viral marketing

- Example of Viral Marketing: Hotmail.com
 - Jul 1996: Hotmail.com started service
 - Aug 1996: 20K subscribers
 - Dec 1996: 100K
 - Jan 1997: 1 million
 - Jul 1998: 12 million

Bought by Microsoft for \$400 million

Marketing: At the end of each email sent there was
a message to subscribe to Hotmail.com
“Get your free email at Hotmail”

Models of influence in networks

Intuition: fraction of friends that have already adopted the product influence the likelihood of a node becoming an adopter.

Problem:

Select an initial group of k influencers so that - given some propagation model - the expected number of converts is maximized.

Granovetter: Threshold Models for Collective Behavior (1978)

Domingos, Richardson: Mining the Network value of Customers (2001) Mining Knowledge-sharing Sites for Viral Marketing

Kempe, Kleinberg, Tardos: Maximizing the Spread of Influence Through a Social Network (2003)

Models of influence in networks

Intuition: fraction of friends that have already adopted the product influence the likelihood of a node becoming an adopter.

Models:

- Linear Threshold Model
- Independent Cascade Model



Problem:

Select an initial group of k influencers so that - assuming one of the above propagation models - the expected number of converts is maximized.

Linear threshold model

Setup: Given a graph $G(V, E)$

- there is an initial set of active nodes (called seeds)
- once a node becomes active it will possibly activate its neighbors

Linear threshold model

- each node v has an *activation threshold* $\theta_v \in [0,1]$
- node v is influenced by each neighbor w by some weight $b_{v,w}$ such that

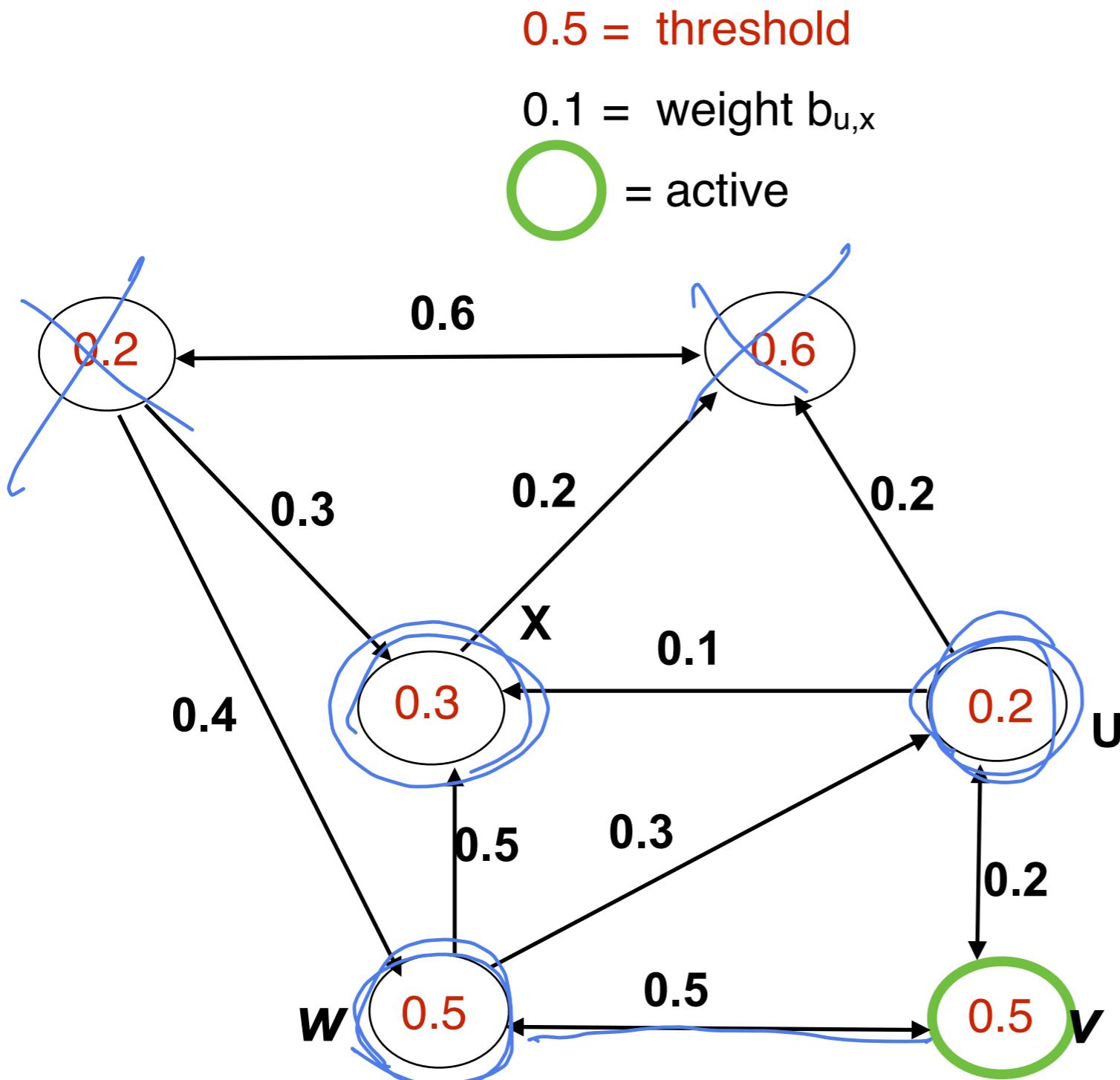
$$\sum_{\substack{w \text{ is neighbor of } v}} b_{v,w} \leq 1$$

- v becomes active iff

$$\sum_{\substack{w \text{ is active} \\ w \text{ is neighbor of } v}} b_{v,w} \geq \theta_v$$

input: $G(V, E)$, θ_v , $b_{v,w}$

Example



Independent cascade model

Setup: Given a graph $G(V, E)$

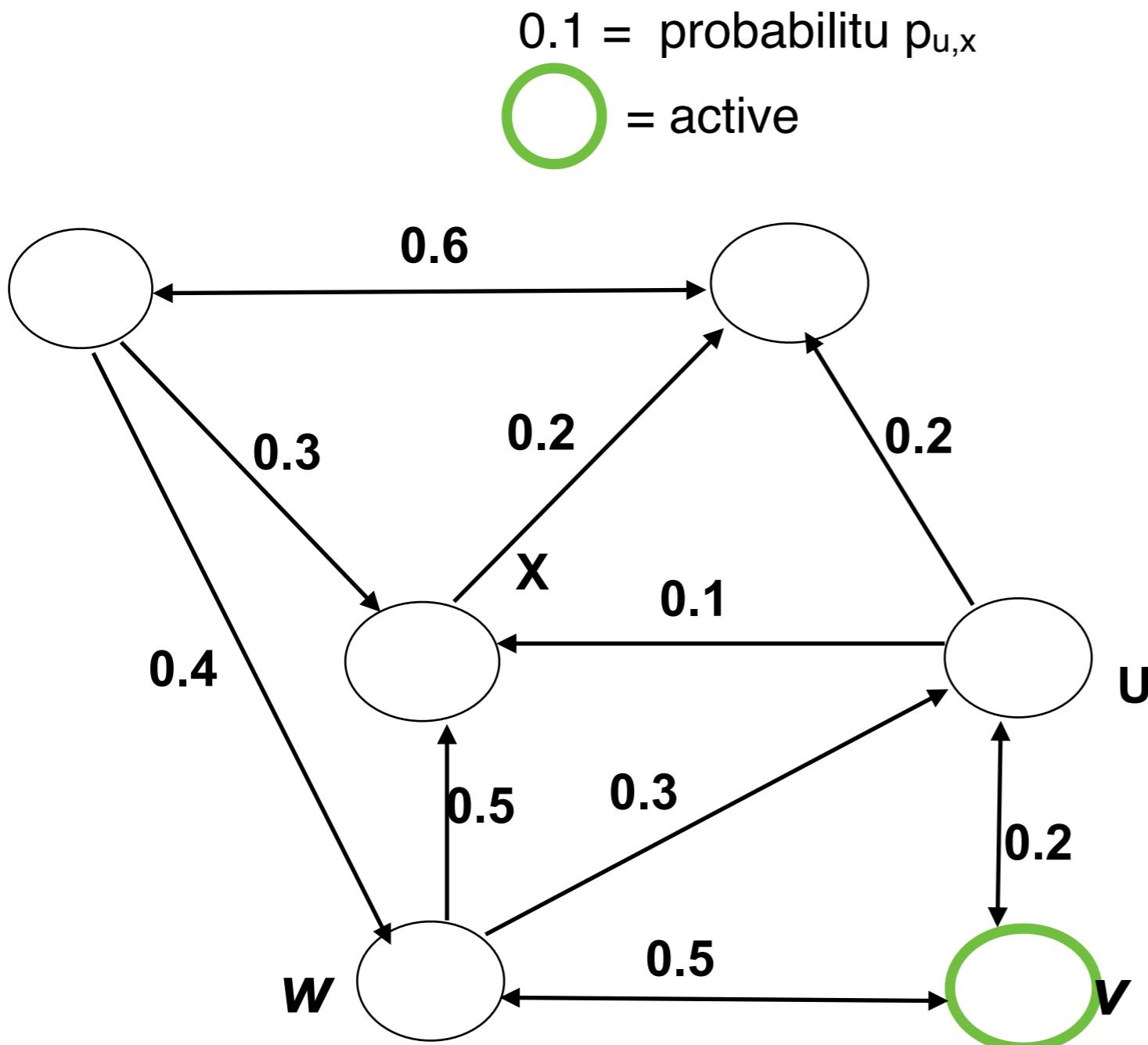
- there is an initial set of active nodes (called seeds)
- once a node becomes active it will possibly activate its neighbors

Independent cascade models

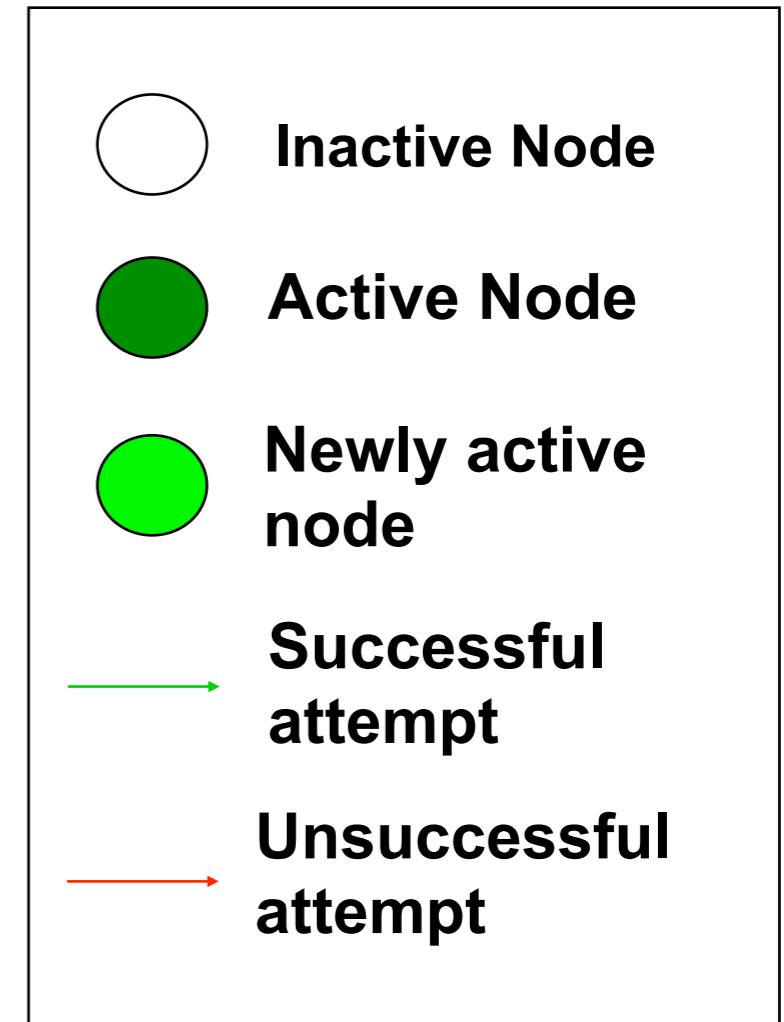
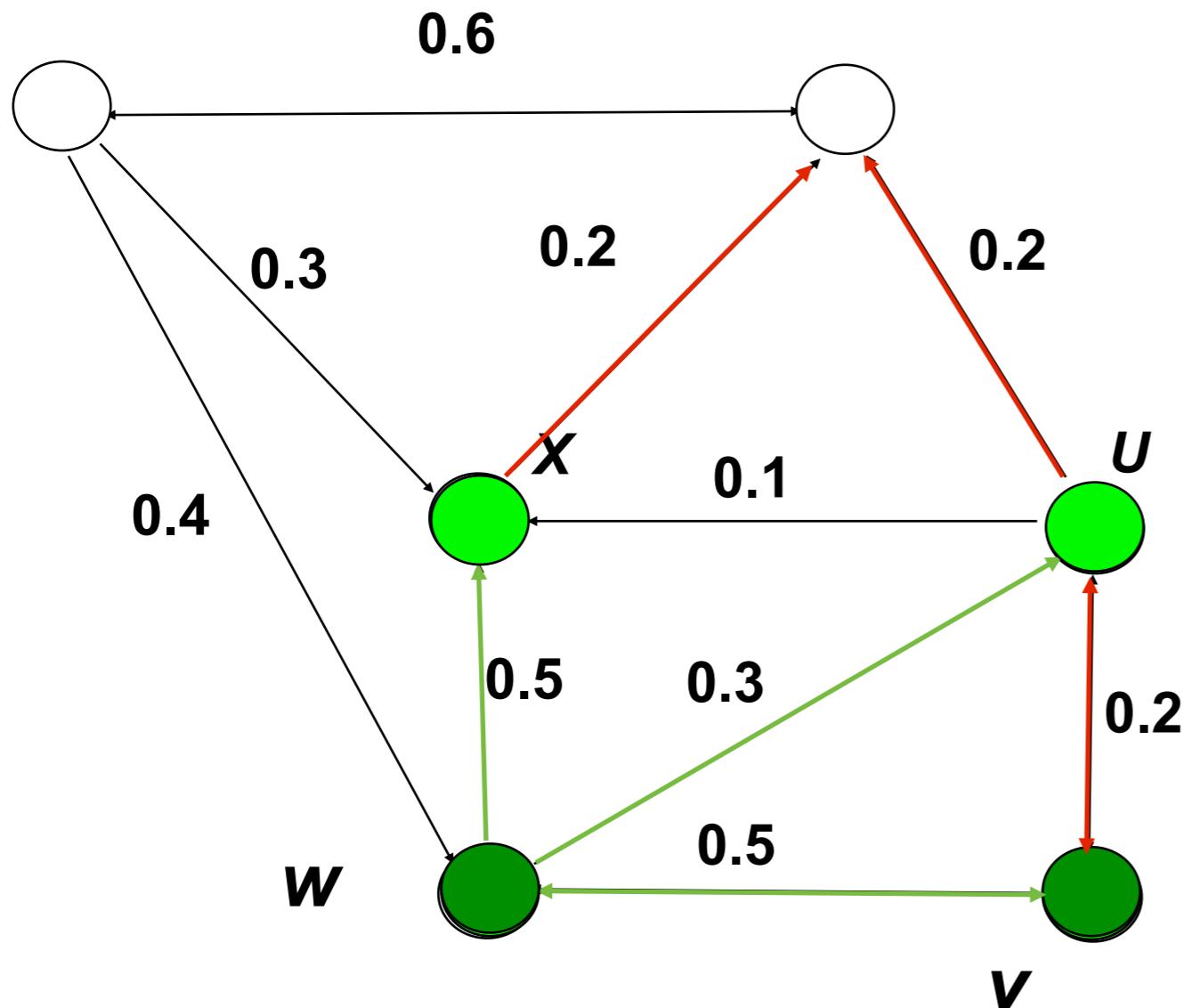
- when a node v becomes active at time t it has a *single* chance to activate its neighbor w
- the activation succeeds with probability $\underline{p_{v,w}}$
- note: if w has multiple active neighbors, each attempts to activate w independent of each other

input: $G(V, E)$, $p_{v,w}$

Example



Example



Stop!

nodes reachable from
the initially active node
along green edges are
active at the end

Influence maximization problem

Let $G(V, E)$ be a graph

the **influence** $f(S)$ of node set S is the *expected* number of active nodes, given one of the two models, if S is the initially active set.

Influence maximization problem: Given as input $G(V, E)$, one of the models with parameters and budget k , find a set S of k nodes with maximum influence $f(S)$

What can we say about the objective function $f(S)$?

↳ both methods
have submodular
objective functions

Influence maximization problem

Let $G(V, E)$ be a graph

the **influence** $f(S)$ of node set S is the *expected* number of active nodes, given one of the two models, if S is the initially active set.

Influence maximization problem: Given as input $G(V, E)$, one of the models with parameters and budget k , find a set S of k nodes with maximum influence $f(S)$

What can we say about the objective function $f(S)$?

- **monotone increasing** — adding one more node to S can only increase the influence
- **submodular** — adding an additional node to a smaller set S has larger impact on the spread
 - how can we prove this given that $f(S)$ is a probabilistic function? → use expected value?
$$f(S \cup \{v\}) - f(S) \geq f(T \cup \{v\}) - f(T) \text{ for } S \subset T \subset V$$

Greedy optimization

Suppose we can prove that the probabilistic function $f(S)$ is submodular

Greedy algorithm:

- 1.start $S = \emptyset$
- 2.for k iteration add v such that in expectation $f(S \cup \{v\}) - f(S)$ is max

Theorem: this greedy algorithm yields a $(1-1/e)$ -approximation

The expected number of activate nodes, when the seeds are selected with the greedy algorithm are ~63% of the expected number for the best seed set.

Proof of submodularity for random independent cascade model

cascade process: if a node v is activated, then flip a coin for each adjacent edge (v,w) to activate w with probability $p_{v,w}$

instead, generate “possible world” G_r

- iterate over each edge of G first
- for each (v,w) flip a coin and keep the edge with probability $p_{v,w}$
- now we have a deterministic graph - an instance of the random graph

active nodes at the end of the diffusion are the ones *reachable* from the seeds in this generated graph

- reachability is submodular - the seeds are nodes that “cover” the paths

conclusion: for any one specific instance of the random model, the influence function $f(S)$ is submodular.

Proof of submodularity for random independent cascade model

cascade process: if a node v is activated, then flip a coin for each adjacent edge (v,w) to activate w with probability $p_{v,w}$

instead, generate “possible world” G_r

- iterate over each edge of G first
- for each (v,w) flip a coin and keep the edge with probability $p_{v,w}$
- now we have a deterministic graph - an instance of the random graph

method to simulate the process for experiments!

active nodes at the end of the diffusion are the ones *reachable* from the seeds in this generated graph

- reachability is submodular - the seeds are nodes that “cover” the paths

conclusion: for any one specific instance of the random model, the influence function $f(S)$ is submodular.

Proof of submodularity for random independent cascade model

conclusion: for any one specific instance of the random model, the influence function $f(S)$ is submodular.

fact: non-negative linear combination of submodular functions is also submodular

expected influence in G :

S = set of seed nodes

G_r = random graph instance

$A(G_r)$ = set of active nodes in G_r given S as seed set

$$f(S) = \sum_{G_r} Pr(G_r) \cdot |A(G_r)|$$

Similar proof can be done for the linear threshold model

Implementing greedy

In practice, how can we implement an optimization algorithm with a random objective function?

Still an open question how to compute efficiently

- Kempe et al.: neat trick called “lazy greedy updates” → only update coverage computation for top few candidate

We get very good estimations by simulation

- repeat many times:
 - generate G_r
 - find the optimal set S on G_r using the deterministic (set cover-style) greedy algorithm
- influence can be computed as the average activation over the many runs

Experimental results - Kempe, Kleinberg, Tardos [2003]

Data:

co-authorship graph in papers on arXiv in the high-energy physics theory section

graph $G(V, E)$

V : authors

E : there is an edge (v, w) if persons v and w have written a paper together

$|V| = 10748$, $|E| = 53000$

model parameters

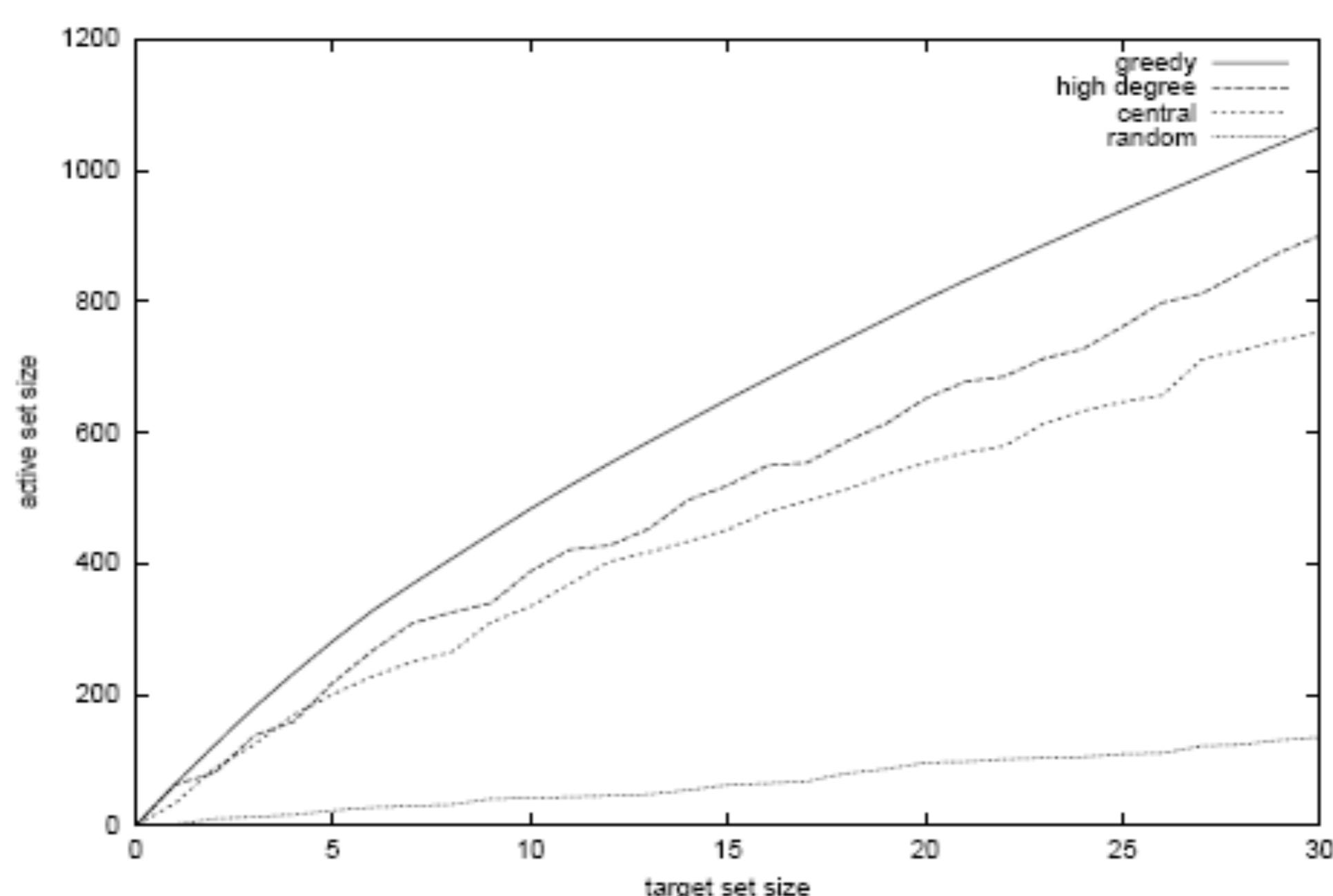
- *linear threshold*: based on multiplicity of edges
 - fraction of papers co-authored $c_{v,w}$ divided by all papers by this person d_v

$$b_{v,w} = \frac{c_{v,w}}{d_v}$$

- *independent cascade*: activation probabilities chosen uniform at random

Experimental results - Kempe, Kleinberg, Tardos [2003]

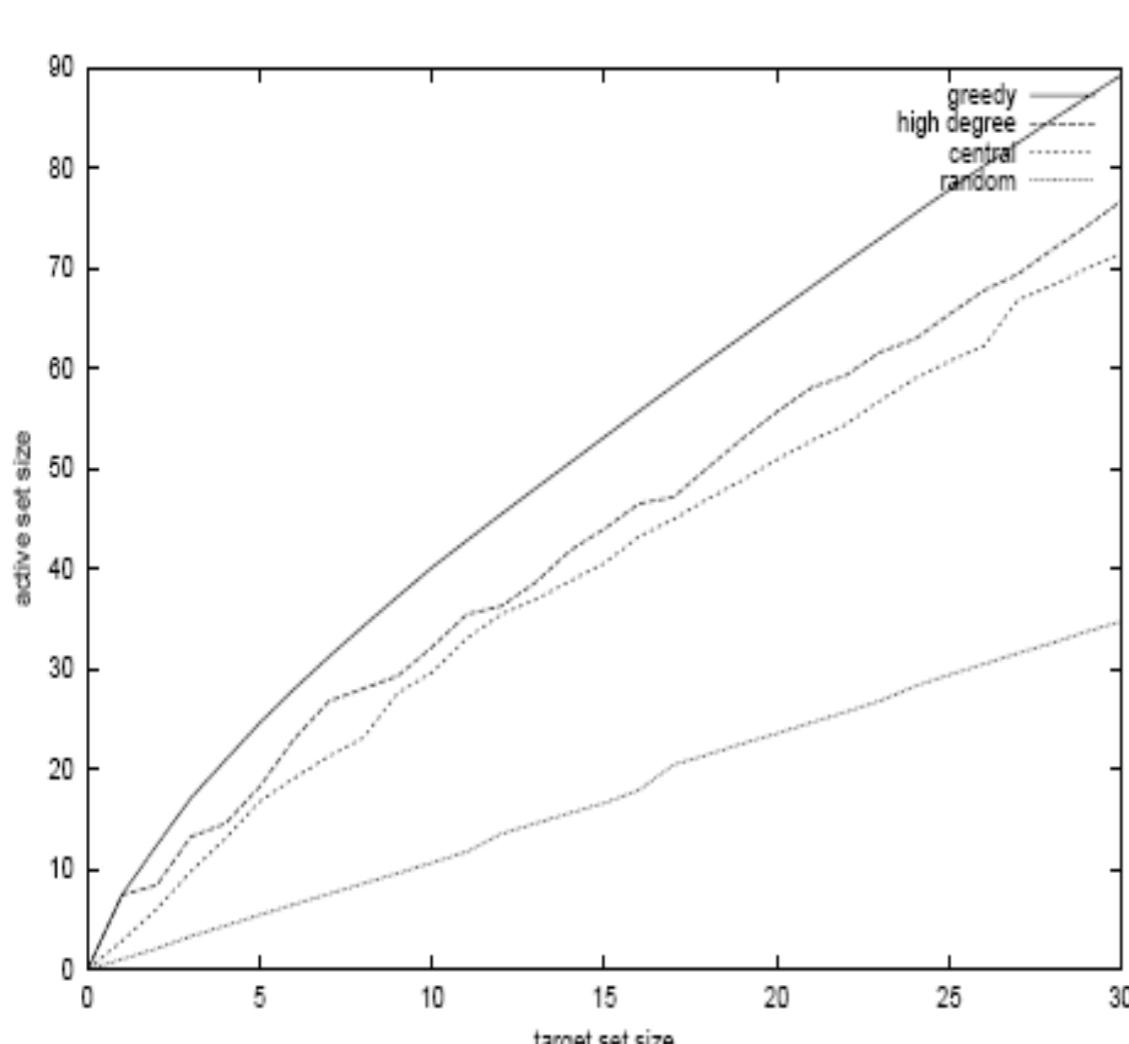
- Simulate process 1000 times, re-select probabilities and thresholds each time
- compare to 3 common heuristics



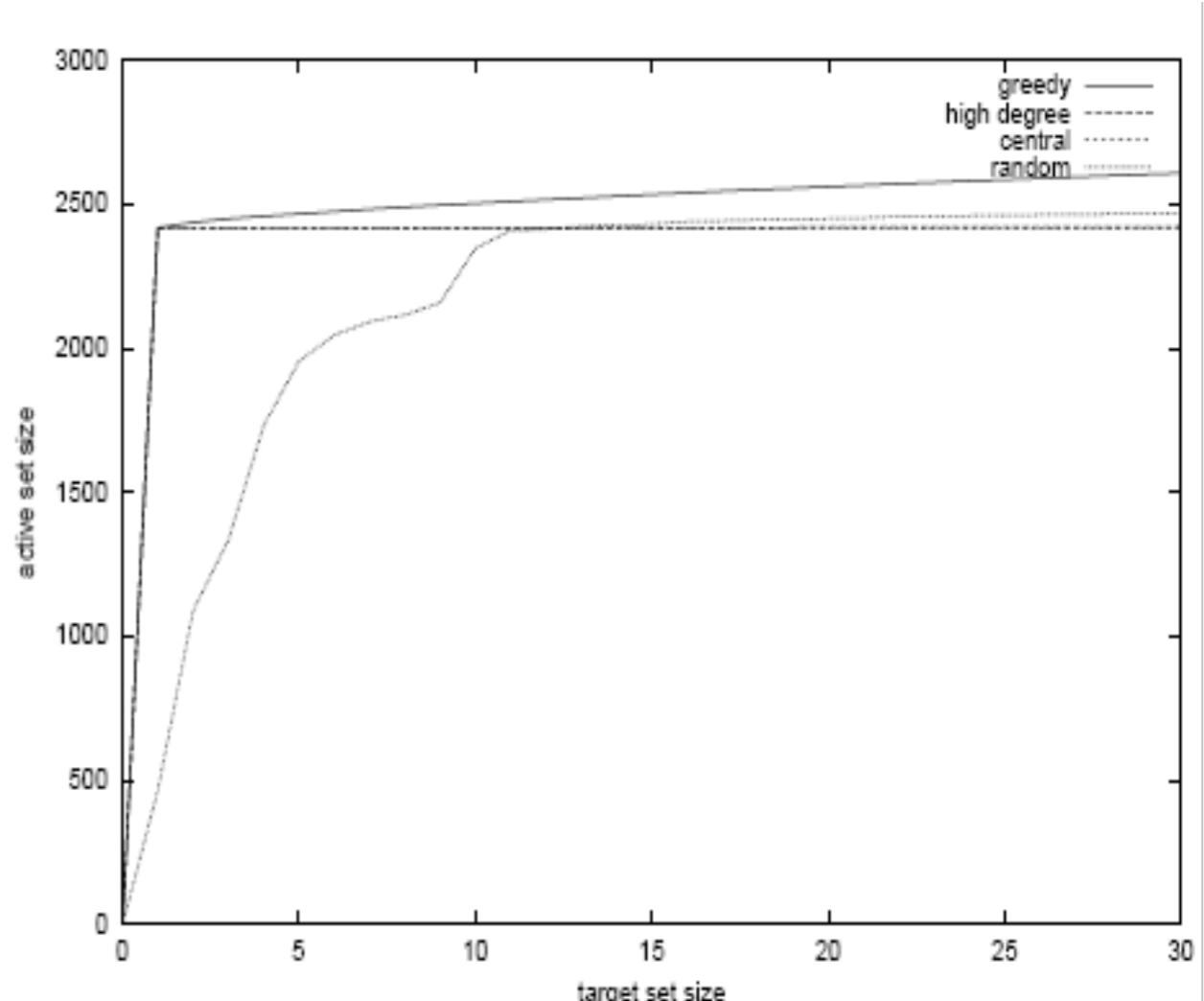
Results for linear threshold model

Experimental results - Kempe, Kleinberg, Tardos [2003]

- Simulate process 1000 times, re-select probabilities and thresholds each time
- compare to 3 common heuristics



$p = 0.01$



$p = 0.1$

Results for independent cascade model