

## Week 9

10/19/2020

- Next thw due Monday 10/26,  
(extended from this Fri 10/23)
- Feedback from proposals → inbox.
- This week: sampled networks.

Why sample?

- ① Access (e.g. private Twitter accounts)
- ② Computation Time (algs are slow or memory runs out)

WWW  $\sim 10^{10}$  vertices

If we have  $G = (V, E)$   
a subsample or subgraph

$$G' = (V', E')$$

where  $V' \subset V$  and.

$$E' \subseteq V' \times V' \subseteq E$$

is a  
 subset  
 with possible  
 equality

If we don't sample  
the nodes we don't get  
any of its edges.

## Key:

Working with  $G'$ , rather than  $G$  can be fine if the function  $f$  you are computing would yield equivalent results on  $G$  and  $G'$ , i.e.

$f(G) = f(G')$ , orrrr... more generally: if the distribution of outputs is the same  $\Pr(f(G)) = \Pr(f(G'))$ .

- Depends on  $f$ .
- Depends on  $G'$ .

Do we have access to the whole network (or can we query it?)

yes

Probabilistic Sampling

no

Seed-based sampling

① uniform vertex.

- include each vertex  $i$  (and neighbors) w.p.  $p$ .

② uniform edge.

- include each edge  $(i,j)$  and endpoints w.p.  $p$ .

③ Degree-proportional. include each vertex : (and its neighbours) w.p.  $\propto k_i$

④ Attribute-proportional. " " " w.p.  $\propto x_i$

## Seed-Based Sampling.

(Assumes access to one or more seed vertices)

$l \leftarrow \text{level}$

### ① Snowball sampling.

For each seed vertex  $i$ , and distance  $l$ , include all vertices and their neighbors, in an  $l$ -step B.F.S. tree, rooted at  $i$ .

### ② B.F.S. edge sampling.

For each seed vertex  $i$ , and distance  $l$ , include all edges in an  $l$ -step BFS, rooted at  $i$ .

### ③ Adaptive sampling:

For each seed  $i$ , and int  $s$ , include all vertices (and nbrs) or include all edges, in an adaptively grow tree containing  $s$  vertices, rooted at  $i$ .

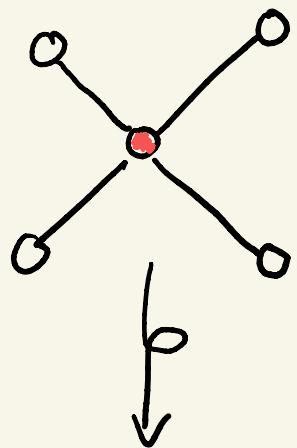
## Analysis

### Sampling induces Patterns.

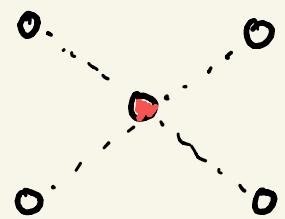
- ① Extreme sparsity (prob. sampling).
- ② Compact, but biased subgraph. (seed-based)
- ③ An overabundance of low-degree vertices in  $G'$ . (often  $k=1$ )

Why?

# Network Immunization!



$n = 5$   
 $v = 1$



$n$  nodes  
 $m$  edges  $\rightarrow$  adj A

$v$  vaccines. 100% effective.

Where to vaccinate?

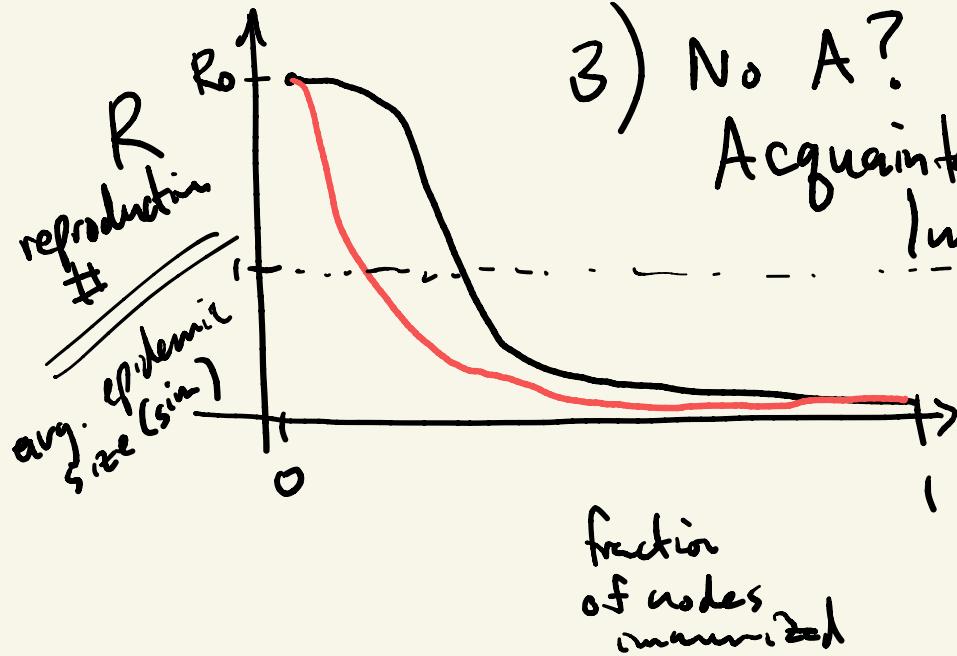
W/H: between countries: allocation.  
within country: prioritization

Strategies?

- 1) High centrality
  - eigenvector
  - degree

- 2) Maximize # of disconnected components.

- 3) No A?  
Acquaintance immunization.



10/21/20

Sampling distorts ...

- degrees - too many  $k=1$  and low-degree nodes
- components. May have more components than in the original network!

## Uniform Vertex Sampling

- target Samp  $s$  vertices
- choose <sup>each</sup> vertex w.p.  $p$ . and, in so doing, also get its nbrs.

How should we choose  $p$ ?

Each time we sample  $i$ , in expectation, we'll add  $1 + \langle k \rangle$  vertices to sample.

$\Rightarrow$  choose each vertex w.p.  $P = \frac{s}{(1 + \langle k \rangle)n}$   
we should get approx  $s$  vertices in the sample.

$$\text{sub } \langle k \rangle = \frac{2m}{n}$$

$$\Rightarrow P = \frac{s}{\left(1 + \frac{2m}{n}\right)n} = \boxed{\frac{s}{n + 2m} = P}$$

If network is sparse,  $P \propto \frac{s}{n}$ .

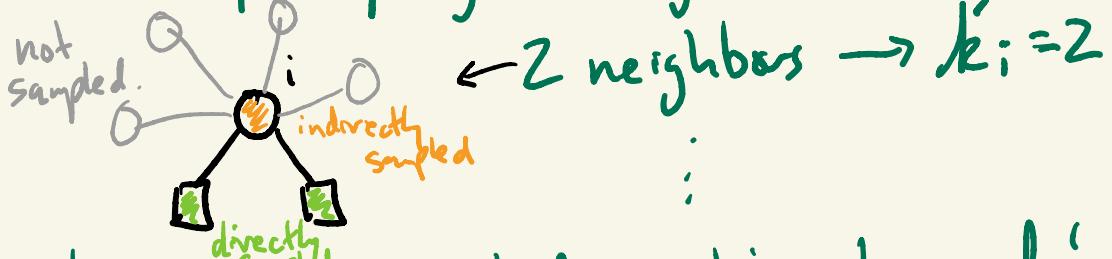
How will sampling change degree dist?

Note:

- directly sampling a vertex  $\Rightarrow$  we get  $k'_i$  right,  $k'_i = k_i$

prime means sampled degree  
no prime = original degree.

- directly sampling 1 neighbor  $\rightarrow k'_i = 1$



How many sampled vertices have  $k' = 1$ ?

Two ways to have  $k' = 1$

① we sampled  $i$ , and it has  $k=1$

$$\begin{matrix} \downarrow \\ \text{P}_{k=1} \end{matrix} \quad Pn_i$$

② we sampled another node, whose neighbor is  $i$ , but only sample one neighbor of  $i$ .

Let  $n_k$  be # of nodes

in original graph w/ degree  $k$ .

$$E[n'_i] = pn_i + \sum_{k=2}^n p(1-p)^{k-1} n_k$$

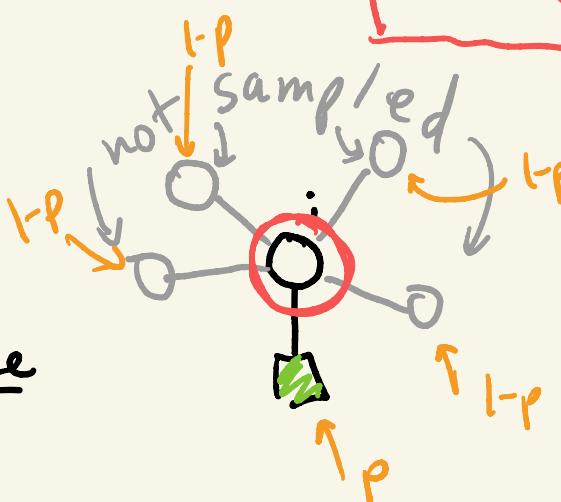
$$E[n'_{k'}] = pn_k + p^k \sum_{K=k'}^n (1-p)^{k-k'} n_k$$

3 pts E.C.

email me, and tell  
me w/ short expl. which

is correct?

where is  $P(\text{did not sample } i)$ ?



$$P(1-p)^{k-1}$$

## Uniform Edge

- choose to include each edge w.p.  $P$ .
- include neighbors



sampling nodes  $\propto k_i$

A node w/ degree  $k$ :

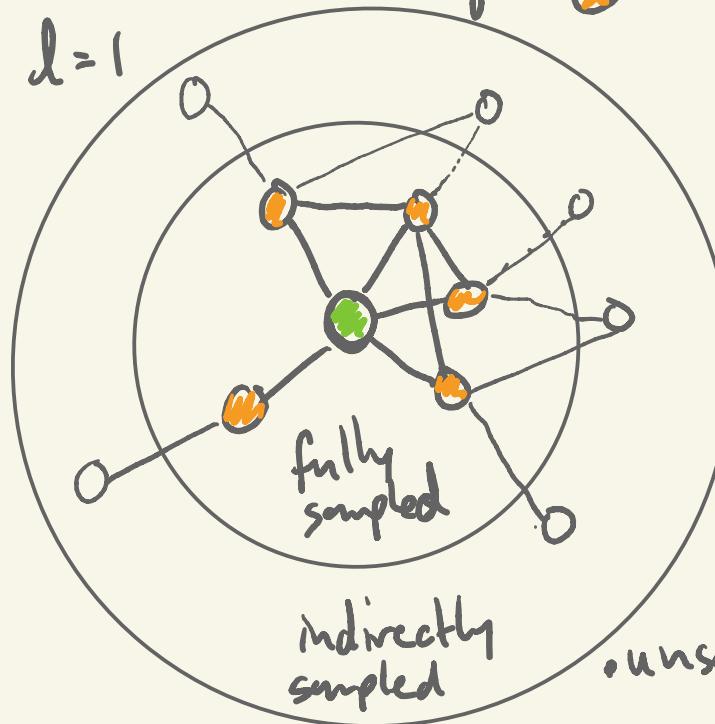
will have, in exp. degree  $k'_i = p k_i$ :  
in sampled graph.

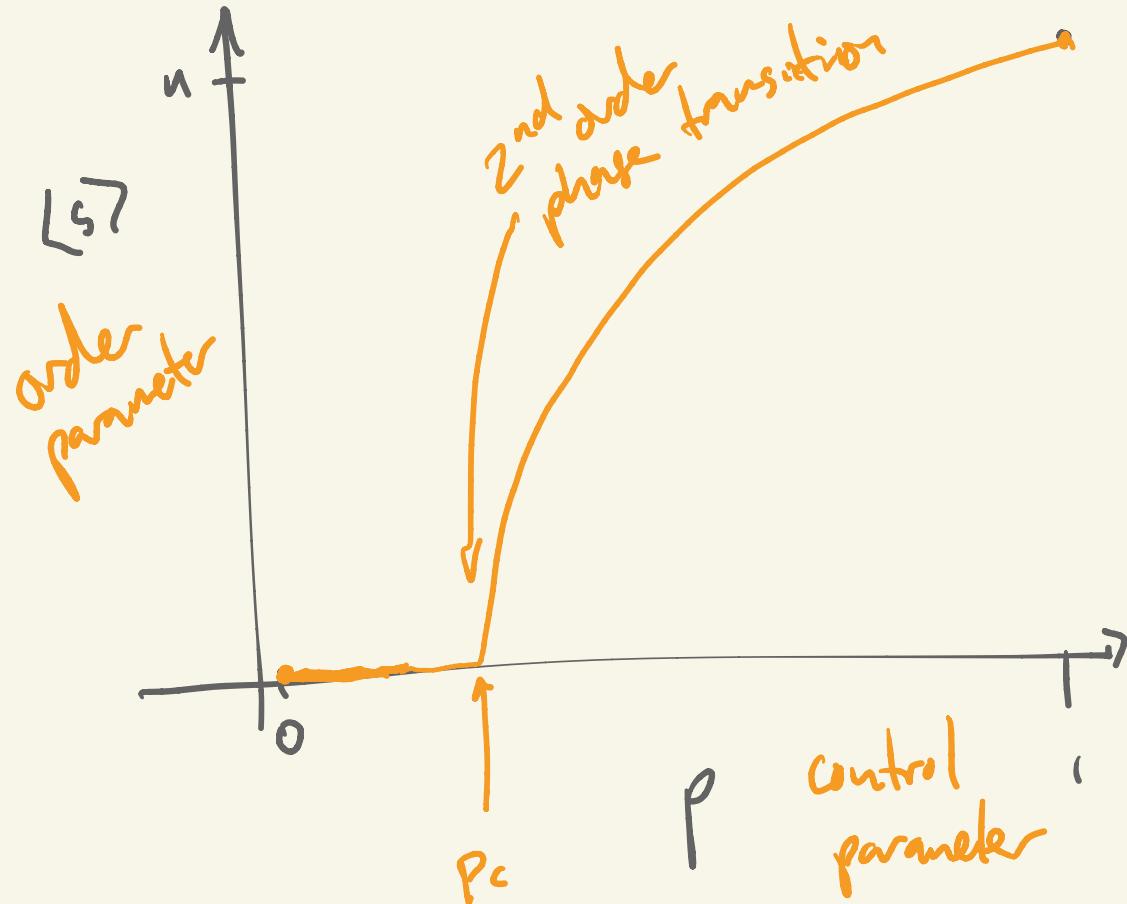
Parallel: percolation.

↳ if  $p$  too small, network will be disconnected, sparse. (if  $\langle k \rangle$  falls below E.R. critical value  $\langle k \rangle_c = 1$ )

## Snowball Sampling!

- start from a seed  $i$  (green)
- include all vertices and nbrs up to a distance  $l$  and edges (orange)
- fully sample: and any  $j$  s.t.  $d_{ij} \leq l$
- indirectly sampled  $j$  s.t.  $d_{ij} = l+1$
- unsampled  $j$  s.t.  $d_{ij} > l+1$





What would  
 $\langle d \rangle$  be for  $p=1$ ?



How many steps to reach  
all nodes in network?

$$\frac{\text{diam}}{2} \leq \langle d \rangle \leq \text{diam}$$

$$\langle d \rangle = \frac{1}{n} \sum_i \text{longest shortest path from } i \text{ to any } j$$