

graph Sampling

SYSM 6302

CLASS 18

Graph Sampling



Given a population graph $G = (V, E)$

Obtain a sample graph $G^* = (V^*, E^*)$

Often we think of G^* as a subgraph
of G , but doesn't have to be
(additional "incorrect" edges)

Consider a statistic of G : $\eta(G)$

Can we use G^* to get a reliable estimate?

$$\eta(G) \stackrel{?}{\approx} \hat{\eta} = \eta(G^*)$$

↑
(the "plug-in" method)

Suppose $\eta(G) = \frac{1}{N} \sum_{i \in V} k_i$ Average Degree $N = |V|$

① Sample n nodes without replacement

$$V^* = \{i_1, i_2, \dots, i_n\} \text{ with degrees } \{k_i^*\}_{i \in V}$$

$$\Rightarrow \eta(G^*) = \frac{1}{n} \sum_{i \in V^*} k_i^* \quad \leftarrow \text{plug-in approach}$$

Consider two sampling methods:

- 2.1 k_i^* is the actual degree of the node (i.e., observe all incident edges)
 - 2.2 k_i^* is the degree of the graph induced by V^* (i.e., only observe edges between sampled nodes)
- 2.2 greatly underestimates $k_i \Rightarrow$ underestimates $\eta(G)$



In 2.2 $k_i^* \approx k_i \cdot \frac{n}{N}$ fraction of edges
expected to fall
into the sample
tactual degree

Adjusting an estimate for sample bias

(Horvitz-Thompson Estimation)
(height, age, income, etc)

Stepping back from networks...

Given a population U with values: $\{y_i\}_{i=1}^N$

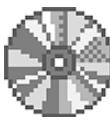
$$\text{Total: } \tau = \sum_{i \in U} y_i \quad \text{Mean: } \mu = \frac{\tau}{N}$$

Given a sample S with values: $\{y_i\}_{i=1}^n$ simple random sample w/ replacement

$$\text{Sample mean: } \bar{y} = \frac{1}{n} \sum_{i \in S} y_i \rightarrow E[\bar{y}] = \mu \text{ and } E[n\bar{y}] = \tau$$

expectation \Rightarrow Unbiased Estimate

If a different sampling technique is used could lead to biased estimates



$$\hat{\bar{Y}}_{\pi} = \sum_{i \in S} \frac{y_i}{\pi_i}$$

← unbiat the estimate through a weighted average
(weighted by the probability π_i)

↑ probability that
 i is in the sample ← inclusion probability
($\pi_i > 0$)

Define $z_i \in \{0, 1\}$, $z_i = \begin{cases} 1 & \text{if } i \text{ is in sample } S \\ 0 & \text{otherwise} \end{cases}$

$$E[z_i] = P(z_i=1) = \pi_i$$

$$E[\hat{\bar{Y}}_{\pi}] = E\left[\sum_{i \in S} \frac{y_i}{\pi_i}\right] = E\left[\sum_{i \in U} \frac{y_i}{\pi_i} z_i\right] = \sum_{i \in U} \frac{y_i}{\pi_i} \underbrace{E[z_i]}_{\text{red}} = \sum_{i \in U} y_i$$

Take-away: if we can quantify π_i , we can unbiat our estimates
↑ depends on sampling method

Induced Subgraph Sampling



n nodes sampled without replacement $\Rightarrow V^*$

$$E^* = \{(u,v) \mid u,v \in V^*, (u,v) \in E\} \quad \leftarrow \text{Induced edges}$$

Nodes selected uniformly : $\pi_i = \frac{\binom{N-1}{n-1}}{\binom{N}{n}}$

of ways that node i can be in the sample

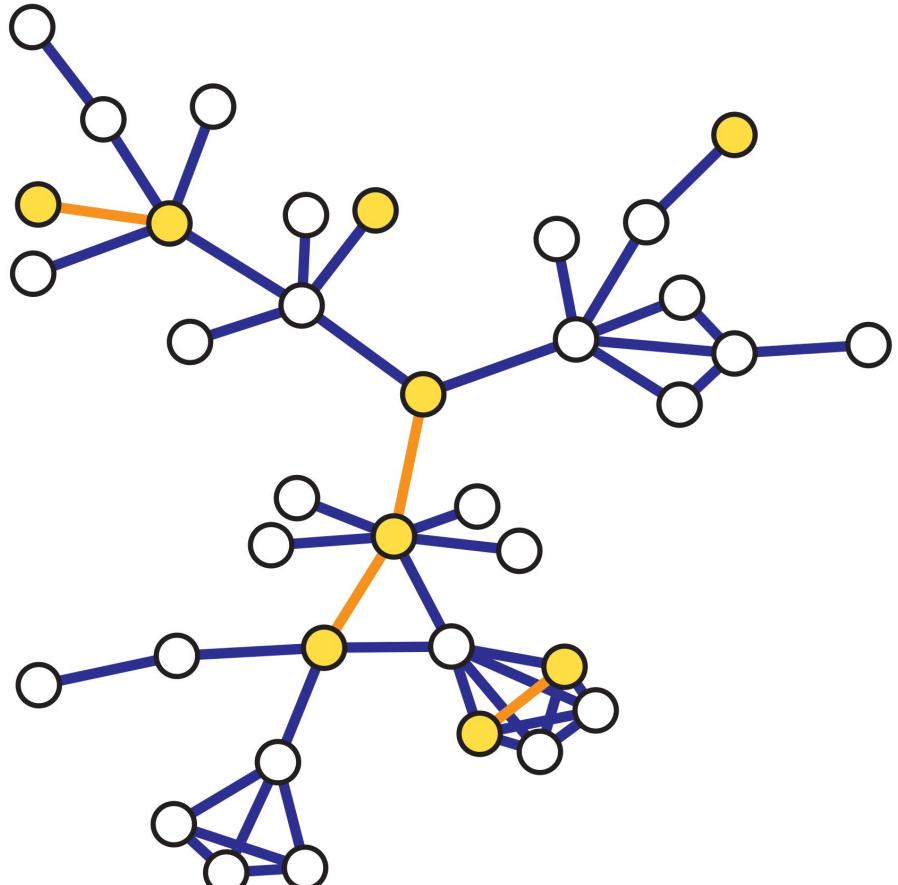
possible samples of size n

$$= \frac{\frac{(N-1)!}{(n-1)! (N-n)!}}{\frac{N!}{n! (N-n)!}} = \frac{n}{N}$$

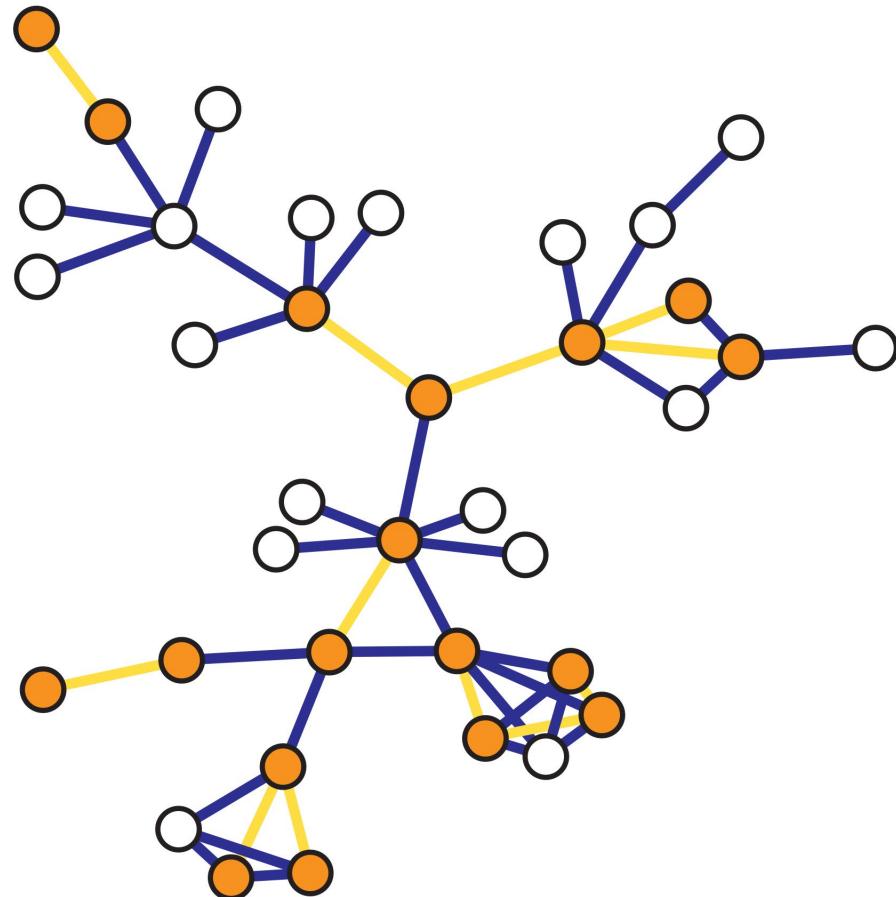
For node i and node j to both be sampled : $\pi_{ij} = \underbrace{\frac{n}{N}}_{\substack{\text{prob. of} \\ \text{one node}}} \cdot \underbrace{\frac{n-1}{N-1}}_{\substack{\text{prob. of} \\ \text{next node}}} = \frac{n(n-1)}{N(N-1)}$

probability ↑
(i,j) in sample

These do require that we know N



Induced



Incident

Incident Subgraph Sampling



n edges are sampled without replacement $\rightarrow E^*$

All nodes incident to sampled edges are observed $\rightarrow V^*$

$$\pi_{i,j} = \frac{n}{|E|}$$

edges sampled
out of so many total edges

$$= \frac{n}{M}$$

Requires knowledge of M

$$\pi_i = P(\text{vertex } i \text{ is sampled}) = 1 - P(\text{no edge incident to node } i \text{ is sampled})$$

$$= \begin{cases} 1 - \frac{\binom{M-k_i}{n}}{\binom{M}{n}} & \text{if } n \leq M-k_i \\ 1 & \text{if } n > M-k_i \end{cases}$$

Vertex inclusion probabilities are non-uniform, depending on its degree

requires knowledge of k_i

Snowball Sampling



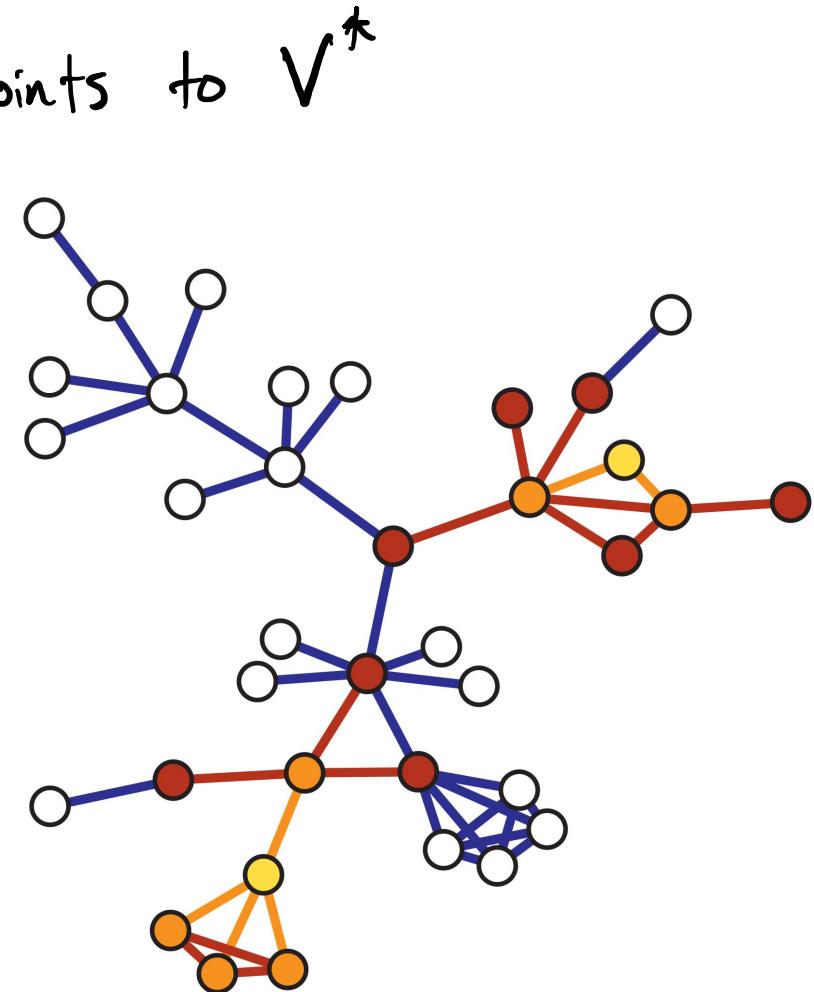
① An initial set of nodes V_0 is selected

② Set $V^* = V_0$

③ Follow edges incident to V^* and add other endpoints to V^*

④ Repeat ③ until some stopping condition is met

This is the process used by web-crawlers



Link Tracing



- After the selection of the initial sample a subset of edges incident to the sample are followed to additional nodes. ↑ according to a predesignated rule
- Snow ball sampling is a special case in which all edges are followed