

# Counting Motifs with Graph Sampling

Jason M. Klusowski and Yihong Wu

Yale University, Department of Statistics and Data Science

## Objectives

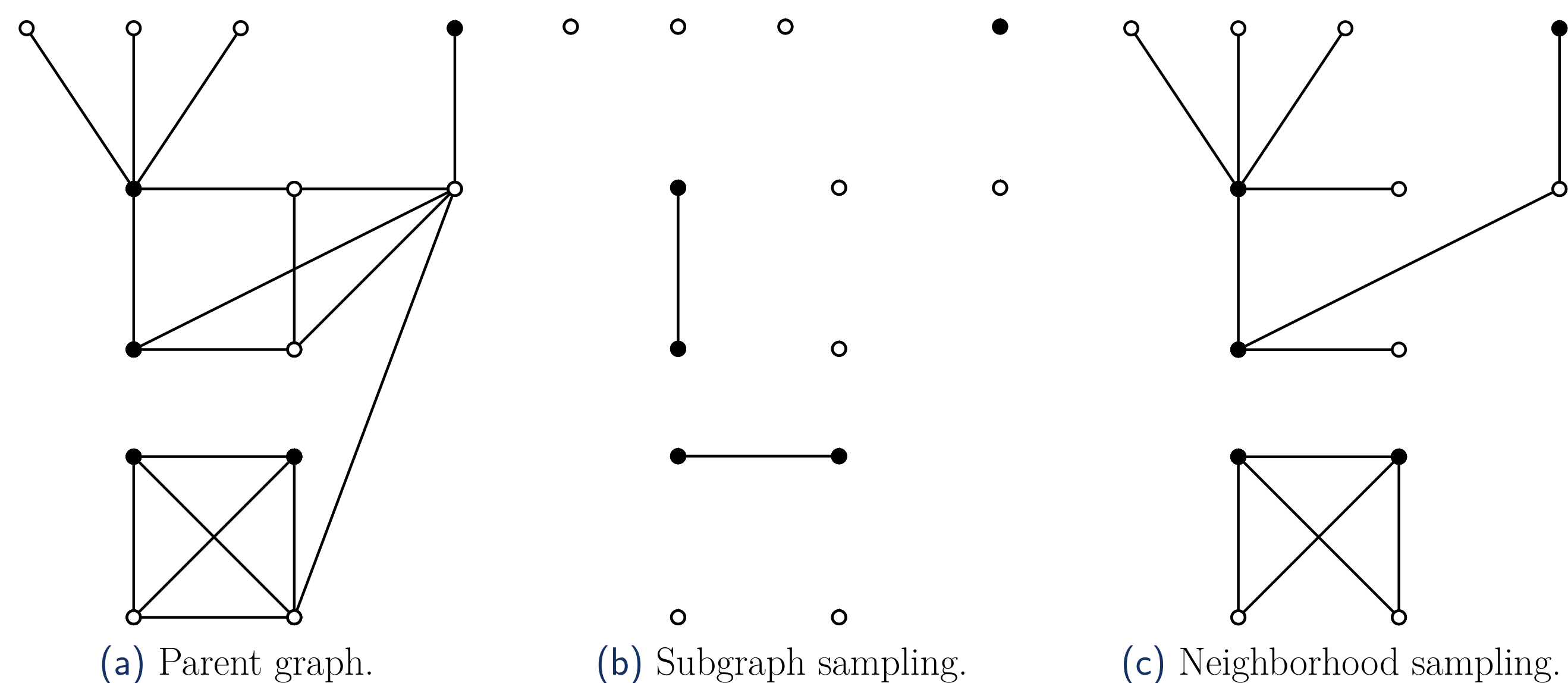
Develop a statistical theory for estimating motif counts (as induced subgraph) in a sampled graph. Focus on large graphs and sublinear sampling regime, where only a vanishing fraction of vertices are sampled.

- How does the sample complexity depend on the motif itself? For example, is estimating the count of open triangles as easy as estimating the closed triangles?
- How much of the graph must be observed to ensure accurate estimation?
- How much more informative is neighborhood sampling than subgraph sampling from the perspective of reducing the sample complexity?

## Introduction

Fix a simple graph  $G = (V, E)$  on  $v(G)$  vertices. We study two sampling models:

- Subgraph sampling.** Sample vertices  $S \subset V$  independently with probability  $p$ . Observe induced subgraph  $G^* \triangleq G[S]$ .
- Neighborhood sampling.** Observe  $G[S]$  and edges between vertices in  $S$  and their neighbors. Sampled graph  $G^*$  is bicolored, with black (sampled) and white (unsampled) vertices.



## Subgraph Counts

Two fundamental quantities govern the distribution of the sampled graph  $G^*$ .

- Subgraph sampling.** Let  $s(g, G)$  denote the number of induced subgraphs of  $G$  that are isomorphic to  $g$ , e.g.,  $s(\text{path}_2, G) = 2$ . Then

$$\mathbb{P}[G^* \simeq g] = s(g, G)p^{v(g)}(1-p)^{v(G)-v(g)},$$

where  $v(g)$  is number of vertices in  $g$ .

- Neighborhood sampling.** Let  $\mathbf{N}(h, G)$  denote the number of ways that a bicolored  $g$  appears (isomorphic as a vertex-colored graph) in  $G$ , e.g.,  $\mathbf{N}(\text{path}_2, G) = 2$ . Then

$$\mathbb{P}[G^* \simeq g] = \mathbf{N}(g, G)p^{v_b(g)}(1-p)^{v(G)-v_b(g)},$$

where  $v_b(g)$  is number of black vertices in  $g$ .

## Optimal Estimators: Achievability

- For parent graph  $G$  with maximal degree bounded by  $d$ , for any motif  $h$  (**connected subgraph**) on  $k$  vertices, estimate  $s = \mathbf{s}(h, G)$  with a multiplicative error of  $\epsilon$ .

- Subgraph sampling.** Horvitz-Thompson (HT) type estimator

$$\hat{s}_h \triangleq \mathbf{s}(h, G^*)/p^{v(h)}.$$

- Neighborhood sampling.** Use HT when  $p \leq 1/d$ . However, HT is suboptimal for  $p > 1/d$ . Tailored estimator is a linear combination of the  $\mathbf{N}(\cdot, G)$  and improves the HT estimator by incorporating the colors of the vertices to reduce (or eliminate) correlation.

- Example.** Edge count estimator has form

$$\hat{e} = \alpha \cdot \mathbf{N}(\bullet-\circ, G^*) + \beta \cdot \mathbf{N}(\bullet-\bullet, G^*),$$

with  $\alpha = \frac{1+dp}{p(2+(d-1)p)}$  and  $\beta = \frac{1-d(1-2p)}{p(2+(d-1)p)}$  optimized to reduce variance.

- Adaptive estimator available with similar guarantees – do not need to know  $d$  a priori.

## Minimax Lower Bounds

- Construct random instances of graphs with matching structures of small subgraphs, akin to the method of moment matching.
- Subgraph sampling.** For any connected graph  $h$  with  $k$  vertices, there exists a pair of connected graphs  $H$  and  $H'$ , such that  $\mathbf{s}(h, H) \neq \mathbf{s}(h, H')$  and  $\mathbf{s}(g, H) = \mathbf{s}(g, H')$  for all connected  $g$  with  $v(g) \leq k-1$ .

- Example for  $h = \text{triangle}$  and  $k = 3$ .**

$$H = \text{triangle} \quad H' = \text{square} \quad \Rightarrow \quad \text{TV}(\mathcal{L}(H^*), \mathcal{L}(H'^*)) = O(p^3).$$

- Neighborhood sampling.** For  $k$ -cliques  $h = K_k$ , there exists two graphs  $H$  and  $H'$  such that  $\mathbf{s}(K_k, H) - \mathbf{s}(K_k, H') \geq 1$  and  $\mathbf{N}(g, H) = \mathbf{N}(g, H')$  for all neighborhood subgraphs  $g$  such that  $v_b(g) \leq k-2$ .

- Example for  $h = \text{triangle}$  and  $k = 3$ .**

$$H = \text{triangle} \quad H' = \text{square} \quad \Rightarrow \quad \text{TV}(\mathcal{L}(H^*), \mathcal{L}(H'^*)) = O(p^2).$$

- Produce two graphs  $H$  and  $H'$  that are statistically indistinguishable unless at least  $k$  vertices (resp.  $k-1$ ) are sampled with subgraph (resp. neighborhood) sampling, but have large separation in the number of motif  $h$ .

## Sample Complexity

- For subgraph sampling, the optimal sampling ratio  $p$  is  $\Theta_k(\max\{(se^2)^{-\frac{1}{k}}, \frac{d^{k-1}}{se^2}\})$ , which only depends on the size of the motif but *not* its actual topology. Furthermore, HT type estimators are universally optimal for any connected motifs.

- When  $s = \mathbf{e}(G)$  is edge count, optimal sampling ratio scales as

$$\Theta\left(\max\left\{\frac{1}{\sqrt{se}}, \frac{d}{se^2}\right\}\right).$$

- For neighborhood sampling, we achieve the sampling ratio  $O_k(\min\{(\frac{d}{se^2})^{\frac{1}{k-1}}, \sqrt{\frac{d^{k-2}}{se^2}}\})$ , which again only depends on the size of  $h$ . Optimal for all motifs with at most 4 vertices and cliques of all sizes.

- When  $s = \mathbf{e}(G)$  is edge count, optimal sampling ratio scales as

$$\Theta\left(\min\left\{\frac{1}{\sqrt{se}}, \frac{d}{se^2}\right\}\right).$$

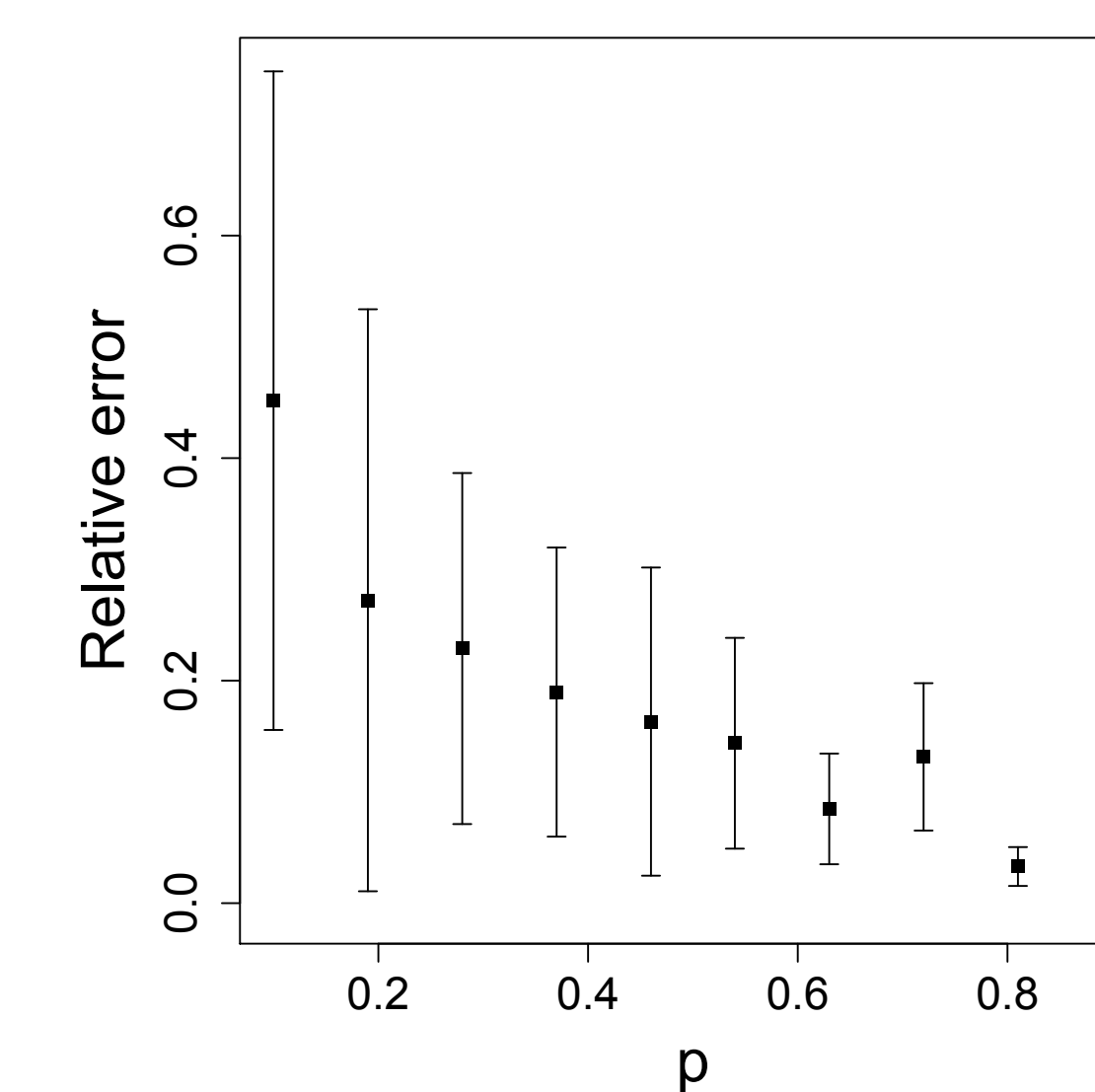
## Additional Structure

To what extent does additional structures of the parent graph, e.g., tree or planarity, impact the sample complexity?

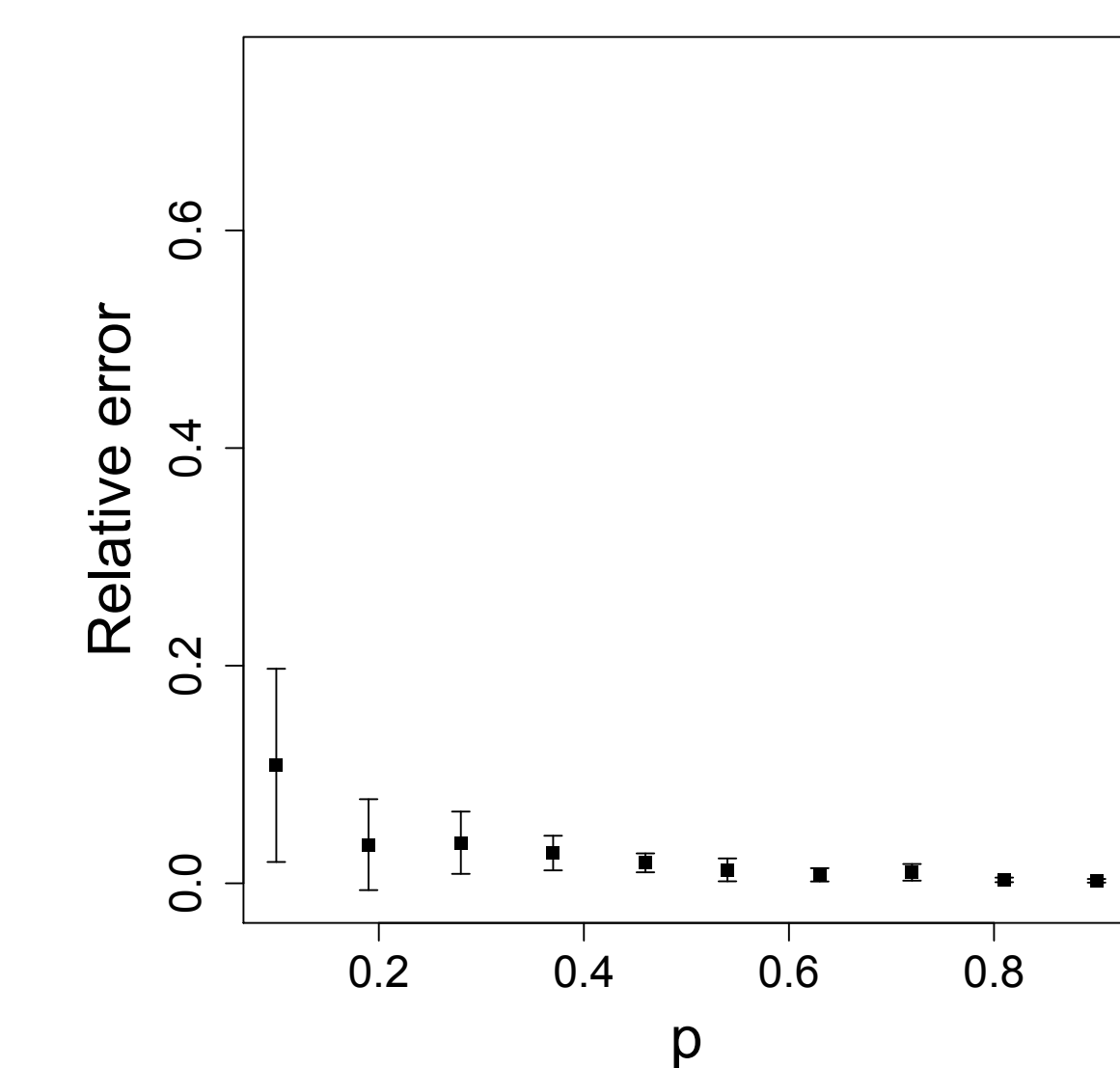
- Tree structure only marginally improves estimation of edges and wedges for both subgraph and neighborhood sampling.
- Planarity improves estimation for triangles for both sampling models.

## Experiments

- Collaboration network between jazz musicians in 198 bands that performed between 1912 and 1940 [Gleiser-Danon, 2003]. Each node is a band and there is an edge between two bands if and only if at least one jazz musician has played in both bands.
- Estimators based on neighborhood sampling perform better than subgraph sampling (significantly less variability).



(a) Jazz network (subgraph sampling).



(b) Jazz network (neighborhood sampling).

Figure: Relative error of estimating the edge count over 10 independent trials. The parent graph  $G$  is the jazz network with  $d = 100$ ,  $v(G) = 198$ ,  $\mathbf{e}(G) = 2742$ .

## Open Questions

- Determine optimal sample complexity in neighborhood sampling for general subgraph counts.
- Statistical limits of  $r$ -hop neighborhood sampling, where we observe a labeled radius- $r$  ball rooted at a randomly chosen vertex [3]. Neighborhood sampling corresponds to  $r = 1$ .
- Statistical limits of counting *edge-induced* subgraphs.

## References

- Jason M. Klusowski and Yihong Wu. Counting motifs with graph sampling. To appear in *Conference on Learning Theory (COLT)*, 2018.
- Jason M. Klusowski and Yihong Wu. Estimating the number of connected components in a graph via subgraph sampling. *arXiv preprint arXiv:1801.04339*, 2018.
- László Lovász. *Large Networks and Graph Limits*, volume 60. American Mathematical Society, 2012.