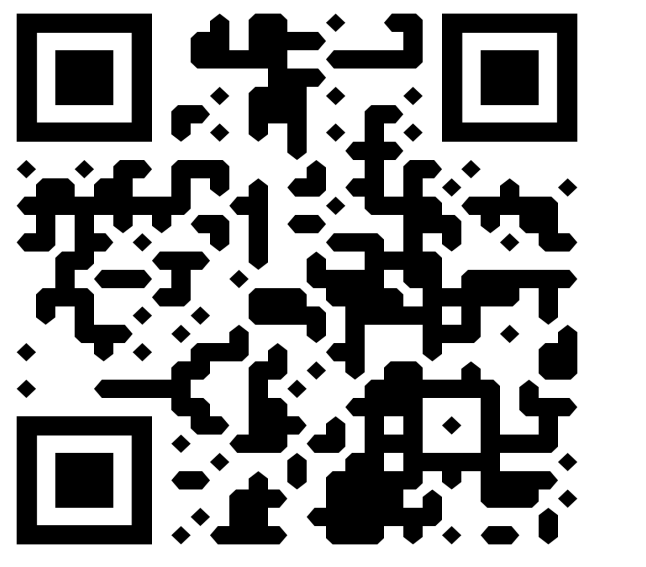


Fast Percolation Centrality Approximation with Importance Sampling

Antonio Cruciani★

Leonardo Pellegrina†



Our Paper

Percolation Centrality is a useful measure to quantify the importance of the vertices in a contagious process or to diffuse information. However, it is **impractical** to compute the exact percolation centrality on modern-sized networks.

Abstract

- There are **key limitations of state-of-the-art** sampling-based approximation algorithms
- We show that, in most cases, the SOTA cannot achieve accurate solutions efficiently

- We propose **PERCIS** a sampling algorithm based on Importance Sampling
- PERCIS severely overperforms the SOTA, both, theoretically and experimentally.

Problem Statement

Input: A graph $G = (V, E)$ with $n = |V|$ and $m = |E|$, and percolation states $\mathbf{x} = (x_1, x_2, \dots, x_n) \in [0, 1]^n$

Problem: Compute the *exact percolation centrality* for each node v ,

$$p(v) = \sum_{s \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}} \cdot \kappa(s, t, v) \in [0, 1]$$

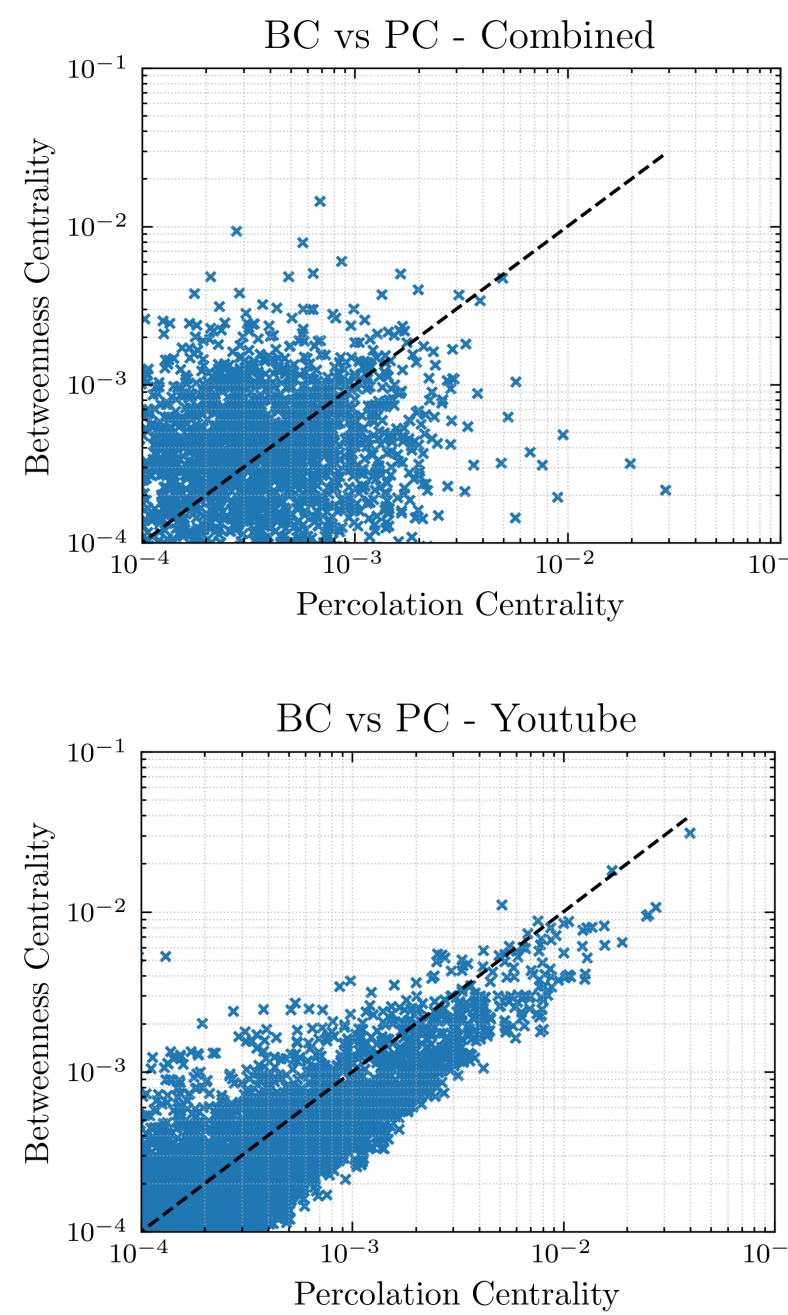
- $\sigma_{st}(v)$ number of shortest paths between s and t passing through v
- σ_{st} overall number of shortest paths between s and t
- $\kappa(s, t, v) = \frac{R(x_s - x_t)}{\sum_{u \neq v} R(x_u - x_v)}$
- $R(x) = \max(0, x)$

Challenge: Exact computation requires $O(n \cdot m)$ time!

Goal: Compute an ε -approximation of the percolation centrality:

$$|p(v) - \tilde{p}(v)| \leq \varepsilon, \quad \forall v \in V$$

Use case: information/contagion propagation in networks



| Graph | Jaccard Similarity Top-K | | |
|----------|--------------------------|-------|-------|
| | 10 | 50 | 100 |
| Guns | 0.053 | 0.087 | 0.117 |
| Combined | 0.0 | 0.031 | 0.015 |
| Youtube | 0.429 | 0.369 | 0.504 |

Jaccard similarity between betweenness and percolation centrality rankings.

State of the art

Lima et al. [1,2] generalised the techniques for the Betweenness centrality to the Percolation centrality.

High level idea:

- Randomly sample shortest paths of the graph
- Use the (weighted) fraction of the paths that traverse v as an estimate of its percolation centrality.

Cons: Technical issues that prevent these methods to be useful in practical applications.

No truly effective algorithm exists to approximate the percolation centrality.

Our Approach: Importance Sampling

Distribution: We define $\tilde{\kappa} : V \times V \rightarrow [0, 1]$

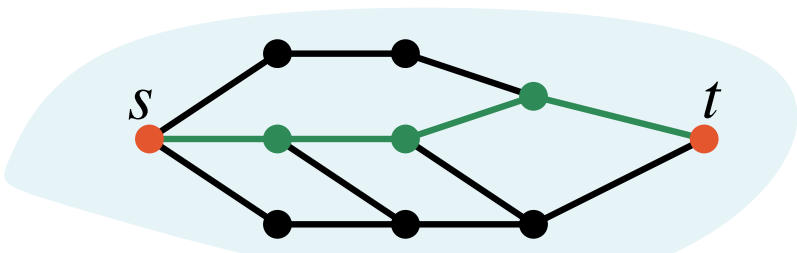
$$\tilde{\kappa}(s, t) = \frac{R(x_s - x_t)}{\sum_{u \neq w} R(x_u - x_w)}$$

For any shortest path τ_{st} , we consider the *importance distribution* :

$$q(\tau_{st}) = \frac{\tilde{\kappa}(s, t)}{\sigma_{st}}$$

Sampling from q

- Sample two nodes s and t with probability $\tilde{\kappa}(s, t)$;
- Compute the set of shortest paths Γ_{st} from s to t ;
- Choose one shortest path uniformly at random from Γ_{st} .



The estimator and its properties

Let $S = \{\tau^1, \tau^2, \dots, \tau^\ell\}$ be a sample of ℓ i.i.d. shortest paths from q .

$$\tilde{p}(v) = \frac{1}{\ell} \sum_{i=1}^{\ell} \frac{\kappa(s, t, v)}{\tilde{\kappa}(s, t)} \mathbb{1}[v \in \mathcal{I}(\tau_{st}^i)]$$

- The estimator is *unbiased*.
- The variance is bounded by $\text{Var}_q[\tilde{p}(v)] \leq \hat{d}p(v)$

Where \hat{d} is the *likelihood ratio*

$$\hat{d} = \max_{v \in V} \left\{ \max_{\substack{s, t \in V, \\ \tilde{\kappa}(s, t) > 0}} \frac{\kappa(s, t, v)}{\tilde{\kappa}(s, t)} \right\}$$

PercIS

Algorithm 1: PERCIS

Input: Graph $G = (V, E)$, percolation states $x_1, x_2, \dots, x_n, \ell_1 \geq 2, \varepsilon, \delta \in (0, 1)$.
Output: ε -approximation of $\{p(v), v \in V\}$ with probability $\geq 1 - \delta$

- $D \leftarrow \text{VERTEXDIAMUB}(G)$;
- $\mathcal{S} \leftarrow \text{IMPORTANCESAMPLER}(G, \{x_v\}, \ell_1)$;
- forall** $v \in V$ **do** $\tilde{p}(v) \leftarrow \frac{1}{\ell} \sum_{i=1}^{\ell} \frac{\kappa(s, t, v)}{\tilde{\kappa}(s, t)} \mathbb{1}[v \in \mathcal{I}(\tau_{st}^i)]$
- $\hat{p} \leftarrow \tilde{p}(S) + \sqrt{\frac{2\Lambda(S) \log(8/\delta)}{\ell_1} + \frac{7D \log(8/\delta)}{3(\ell_1 - 1)}}$;
- $\hat{v} \leftarrow \hat{d}^2 \max_{v \in V} \left\{ \tilde{p}(v) + \sqrt{\frac{2\tilde{p}(v) \log(4/\delta)}{\ell_1} + \frac{\log(4/\delta)}{3\ell_1}} \right\}$;
- $\hat{x} \leftarrow \hat{d}/2 - \sqrt{\hat{d}^2/4 - \min\{\hat{d}^2/4, \hat{v}\}}$;
- $\ell \leftarrow \sup_{x \in (0, \hat{x}]} \left\{ \frac{\hat{d}^2 \ln(\frac{4\hat{d}}{x\delta})}{g(x)h(\frac{x}{\delta\ell_1})} \right\}$;
- $\mathcal{S} \leftarrow \text{IMPORTANCESAMPLER}(G, \{x_v\}, \ell)$;
- forall** $v \in V$ **do** $\tilde{p}(v) \leftarrow \frac{1}{\ell} \sum_{i=1}^{\ell} \frac{\kappa(s, t, v)}{\tilde{\kappa}(s, t)} \mathbb{1}[v \in \mathcal{I}(\tau_{st}^i)]$
- return** $\{\tilde{p}(v), v \in V\}$

Theoretical guarantees of PercIS

IMPORTANCESAMPLER draws ℓ samples from q in time $O(n + \ell(\log n + T_{BFS}))$ and space $O(n + m)$.

Define \hat{v} and \hat{p} such that

$$\max_{v \in V} \text{Var}_q[\tilde{p}(v)] \leq \hat{v}, \quad \sum_{v \in V} p(v) \leq \hat{d}\hat{p} \quad \text{avg. path length!}$$

Given a sample $S = \{\tau^1, \dots, \tau^\ell\}$ of ℓ shortest paths sampled from q , and $\delta, \varepsilon \in (0, 1)$ then

$$\ell \approx \frac{(2\hat{v} + \frac{2}{3}\varepsilon\hat{d})}{\varepsilon^2} (\ln(\hat{d}\hat{p}/\hat{v}) + \ln(2/\delta))$$

gives an ε -approximation of the percolation centrality with probability $\geq 1 - \delta$

PercIS vs UNIF

State Gap:

$$\Delta = \min_{v \in V} \max_{s \neq v \neq t} (x_s - x_t)$$

- When $\Delta \in \Omega(1)$, the likelihood ratio \hat{d} of the IS distribution q is $\hat{d} \in O(1)$
- There exists instances with $\Delta \in \Omega(1)$ where the likelihood ratio of the uniform distribution is $\Omega(n)$
- There exists instances with $\Delta \in \Omega(1)$ where at least $\Omega(n^2)$ random samples are needed by UNIF, while $O(n)$ random samples are sufficient for PERCIS

For all the considered real world networks, it holds $\Delta = 1$

Networks and Experiments

| Graph | $ V $ | $ E $ | D | ρ | Type |
|----------------|--------|---------|-----|--------|------|
| P2P-Gnutella31 | 62586 | 147892 | 31 | 7.199 | D |
| Cit-HepPh | 34546 | 421534 | 49 | 5.901 | D |
| Soc-Epinions | 75879 | 508837 | 16 | 2.755 | D |
| Soc-Slashdot | 82168 | 870161 | 13 | 2.135 | D |
| Web-Notredame | 325729 | 1469679 | 93 | 9.265 | D |
| Web-Google | 875713 | 5105039 | 51 | 9.713 | D |
| Musae-Facebook | 22470 | 170823 | 15 | 2.974 | U |
| Email-Enron | 36692 | 183831 | 13 | 2.025 | U |
| CA-AstroPH | 18771 | 198050 | 14 | 2.194 | U |

- Random Seeds (RS):** small number of nodes with $x_v = 1$ and the rest to 0
- Random Seeds Spread (RSS):** Simulation of infection spreading from random seeds
- Isolated Component (IC):** small isolated component with some nodes $x_v = 1$ and the rest to 0
- Uniform States (UN):** $x_v \sim \text{Uniform}([0, 1])$

