

BAVARIAN: Betweenness Centrality Approximation with Variance-Aware Rademacher Averages

Extended Version

Cyrus Cousins
Dept. of Computer Science
Brown University
Providence, RI, USA
ccousins@cs.brown.edu

Chloe Wohlgemuth
Dept. of Computer Science
Amherst College
Amherst, MA, USA
cwohlgemuth22@amherst.edu

Matteo Riondato
Dept. of Computer Science
Amherst College
Amherst, MA, USA
mriondato@amherst.edu

“[A]llain Gersten, Hopfen, und Wasser” — 1516 Reinheitsgebot

ABSTRACT

We present BAVARIAN, a collection of sampling-based algorithms for approximating the Betweenness Centrality (BC) of all vertices in a graph. Our algorithms use Monte-Carlo Empirical Rademacher Averages (MCERAs), a concept from statistical learning theory, to efficiently compute tight bounds on the maximum deviation of the estimates from the exact values. The MCERAs provide a sample-dependent approximation guarantee much stronger than the state of the art, thanks to its use of variance-aware probabilistic tail bounds. The flexibility of the MCERAs allows us to introduce a unifying framework that can be instantiated with existing sampling-based estimators of BC, thus allowing a fair comparison between them, decoupled from the sample-complexity results with which they were originally introduced. Additionally, we prove novel sample-complexity results showing that, for all estimators, the sample size sufficient to achieve a desired approximation guarantee depends on the vertex-diameter of the graph, an easy-to-bound characteristic quantity. We also show progressive-sampling algorithms and extensions to other centrality measures, such as percolation centrality. Our extensive experimental evaluation of BAVARIAN shows the improvement over the state-of-the art made possible by the MCERAs, and it allows us to assess the different trade-offs between sample size and accuracy guarantee offered by the different estimators.

CCS CONCEPTS

• Mathematics of computing → Probabilistic algorithms; • Information systems → Social networks; • Theory of computation → Shortest paths; Dynamic graph algorithms; Sketching and sampling; Sample complexity and generalization bounds; Approximation algorithms analysis.

KEYWORDS

Concentration bounds, Dynamic graphs, Percolation centrality, Random sampling, Sample complexity, Statistical learning theory

KDD '21, August 14–18, 2021, Virtual Event, Singapore

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore, <https://doi.org/10.1145/3447548.3467354>.

ACM Reference Format:

Cyrus Cousins, Chloe Wohlgemuth, and Matteo Riondato. 2021. BAVARIAN: Betweenness Centrality Approximation with Variance-Aware Rademacher Averages: Extended Version. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, August 14–18, 2021, Virtual Event, Singapore. ACM, New York, NY, USA, 19 pages. <https://doi.org/10.1145/3447548.3467354>

1 INTRODUCTION

A centrality measure [11] assigns to each vertex or edge in a graph a score that quantifies the importance of that vertex/edge. Many measures and algorithms for them have been introduced in the literature, from simple ones like degree to sophisticated ones based on shortest paths, such as closeness centrality [5] and betweenness centrality [2, 29], which is the focus of this work.

Informally, Betweenness Centrality (BC) quantifies the importance of a vertex/edge z as the fraction of all Shortest Paths (SPs) in the graph that go through z (see (1)). It is a measure of *robustness*, in addition to being a centrality measure [10]: if a vertex/edge with high BC is removed from the graph, then the SPs between many pairs of vertices will change, or even cease to exist. BC has many practical applications from community detection [50], to the study of the resilience of the electrical grid [45], to genomics [31].

Computing the exact BC of every vertex/edge in a graph $G = (V, E)$ is quite expensive. Brandes's algorithm (BA) [17], the state of the art, takes time $O(|V||E|)$ if G is unweighted, and $O(|V|^2 \log |E|)$ otherwise. These running times are *too slow to be practical* on even moderately-sized networks. For this reason, soon after the exact algorithm had been introduced, *sampling-based approximation algorithms* were proposed [19, 36]. Approximate BC scores are *sufficient and acceptable* in practice when they come with *strong guarantees* on their quality, i.e., with an upper bound on the *maximum deviation* of any BC estimate from its unknown exact value (i.e., a bound on the maximum estimation error). Good estimates are of particular value when the graph evolves in a *fully-dynamic* way, that is, when edges and vertices are arbitrarily inserted and removed over time. In this setting, computing the exact BC of each vertex/edge, not only is expensive but also has *limited value*, as these scores continuously change. High-quality approximations are sufficient, and are easy to maintain after updates to the graph.

Sampling-based algorithms exhibit an inherent trade-off between the sample size and the quality guarantee: the larger the sample, the smaller the upper bound on the maximum deviation is, i.e., the better

the guaranteed quality of the estimates. Improved characterizations of this trade-off allows for stronger quality guarantees. Previous works studying the trade-off [12, 19, 21, 30, 56, 57] (see Sect. 2 for an in-depth discussion of these and other relevant contributions) also often introduce novel sampling-based *estimators* for BC, and limit the characterization of the trade-off to a single estimator, usually the one they propose. Different estimators draw samples from different populations (e.g., single vertices, pairs of vertices, or single SPs, see Sect. 3.2), so it is not possible to fairly compare the quality guarantees they offer, i.e., the different characterizations of the sample-size-vs-accuracy-guarantee trade-off. While interesting from theoretical and practical points of view, this proliferation of estimators does not help users in choosing which estimator to employ, and which characterization gives the best “bang for the buck”, especially as there may be even better characterizations.

Contributions. We present BAVARIAN, a suite of sampling-based algorithms for approximating the BC of all vertices in a large graph. Our contributions are the following.

- BAVARIAN is an *algorithmic framework* that can be instantiated with BC estimators (Sect. 3.2), enabling the study of their different computational and statistical trade-offs (Sect. 5). This unifying approach can be used for both static and progressive sampling algorithms (Sect. 4.2), and extended to variants of BC, to percolation centrality, and to dynamic graphs (Sect. 4.3).
- Our analysis of BAVARIAN uses Monte-Carlo Empirical Rademacher Averages (MCERAs) (Sect. 3.3), a key concept from statistical learning theory [62]. The variance-aware tail bounds (Thm. 3.1) for the MCERAs exploit the low variance of BC functions (Sect. 3.3), resulting in a tight characterization of the sample-size-vs-accuracy trade-off, and thus in better quality guarantees at smaller sample sizes than previously possible.
- Our unifying framework and the use of the MCERAs allows us to show that, for every estimator, the sample size sufficient to (probabilistically) obtain a desired accuracy depends on the *vertex diameter* (Sect. 4.1), a characteristic and easy-to-bound quantity of the graph. Previous estimator-specific bounds depend on the size of the largest connected component [57], or on the number of vertices [19], which are much larger than the vertex diameter. This result generalizes one by Riondato and Kornaropoulos [56, Coroll. 1] for a specific estimator. At its core, there are *SP-enforcing families*, a novel concept of independent interest. This result is not just of theoretical interest: we use it to develop sample schedules for practical progressive sampling algorithms (Sect. 4.2).
- Our experimental evaluation (Sect. 5) shows the improvement in the quality guarantee offered by BAVARIAN over the state-of-the-art, which is made possible by the use of the MCERAs. We also analyze the behavior of different estimators for BC, in a first fair comparison enabled by our unifying framework.

Some proofs are deferred to App. A.2, and many experimental results are presented in App. A.3.

2 RELATED WORK

We now discuss the contributions on BC most related to ours. BC was originally introduced in sociology [2, 29], and many variants have been developed since then [13, 18, 26, 42, 49, 51, 52]. In this

work, we focus on vertex BC, but the versatility of our approach allows us to tackle other variants (see Sect. 4.3).

The first efficient algorithm for BC [17], known as Brandes’ algorithm (BA), uses an ingenious recursive formulation to obtain the BC of all vertices in time proportional to n Single-Source-Shortest-Paths (SSSP) computations in an n -vertex graph. The time complexity on unweighted graphs is $O(n^2 + nm)$ where m is the number of edges in the graph, and becomes $O(nm + n^2 \log n)$ on weighted graphs. Despite the remarkable efforts to make this algorithm more scalable in practice [27, 55, 58], its cost is still excessive even for non-humongous graphs. Approximation algorithms [8, 9, 12, 19, 21, 30, 36, 48, 56, 57] and heuristics [3, 28, 46, 47] propose to address this shortcoming. We focus here on works that offer *guarantees on their output* (see also Sect. 3.2). We refer the reader to [56, Sect. 2] and [11] for an in-depth discussion of the heuristics.

All approximation algorithms for BC are based on random sampling, but they draw samples from different populations according to different distributions, compute the approximations using different estimators, and their authors use different approaches to analyze the quality guarantees they offer. Our work has two goals: (1) present BAVARIAN, a unifying framework that can be instantiated with all these estimators, thus allowing a fair comparison between them; and (2) use the MCERAs to obtain tighter quality guarantees and better dependency on graph properties for all estimators.

Brandes and Pich [19] and independently Jacob et al. [36] present the first sampling-based approximation algorithm for computing high-quality estimates of the BC of all nodes. Their estimator is closely related to the inner workings of BA, and the analysis of the sample size vs. quality trade-off uses Hoeffding’s bound [35] and the union bound, thus it is quite loose. Geisberger et al. [30] present a slight refinement of the estimator but the sample size is unchanged. Chehreghani et al. [21] present a sampling-based algorithm for estimating the BC of a *single* node. Extending the guarantee to all nodes via union bound results in the same sample size as the algorithm by Jacob et al. [36] and Brandes and Pich [19].

Riondato and Kornaropoulos [56] propose a new estimator (see Sect. 3.2) and an approximation algorithm whose sample size is obtained using an upper bound to the Vapnik-Chervonenkis dimension [63] of the problem. This quantity depends on the vertex-diameter of the graph (the largest number of vertices on a shortest path). The resulting sample size is therefore much smaller than the one derived by previous works, while still being the result of a worst-case analysis. Bergamini and Meyerhenke [8] show how to better approximate the vertex-diameter on directed networks, but they use the same sample size, as does, in the worst case, the progressive-sampling algorithm by Borassi and Natale [12].

Riondato and Upfal [57] introduce another novel estimator (see Sect. 3.2) and an approximation algorithm whose analysis uses Rademacher averages [40], another core concept of statistical learning theory. Rademacher averages more tightly characterizes the trade-off between sample size and approximation quality using only *sample-dependent* quantities, rather than relying on *graph-dependent* (i.e., population-dependent) properties. As a result, the same approximation quality can be obtained with much smaller samples than before. Our work does not use *deterministic* upper bounds

to the Rademacher averages. Rather, we leverage Monte-Carlo estimation (see (5)) and variance-aware tight deviation bounds (see Thm. 3.1), to achieve both aforementioned goals.

Recently, Fan et al. [28] used graph neural networks to approximate the top- k nodes with highest BC, but without guarantees on the output. As argued in Sect. 1, quality guarantees are important, if not necessary, to make approximate solutions acceptable.

Modern graphs often arise from *dynamic processes*, in which vertices and edges are continuously added and/or removed. Many works look at how to compute and update BC, either exactly or approximately, after an update or batch of updates [8, 9, 32, 34, 38, 43, 54, 57, 64]. We discuss how to adapt our results to the fully-dynamic case in Sect. 4.3.

3 PRELIMINARIES

We now state definitions and theorems used throughout the work.

3.1 Betweenness Centrality

Let $G = (V, E)$ be a graph. The edges may be directed or undirected and may have non-negative weights. We denote with $n = |V|$ the number of vertices in G . For ease of discussion, we assume the graph to be readily available so that operations such as sampling a vertex or a pair of vertices uniformly at random is easy. When this assumption does not hold, one can use appropriate methods to draw these samples [22, 23, 39]. We define a path p from a vertex u to a vertex v as an ordered sequence of *distinct* vertices $(u, z_1, \dots, z_{|p|-2}, v)$. Given a path p between vertices u and v , a vertex w is *internal* to p if and only if $w \neq u$, $w \neq v$, and $w \in p$. We denote the set of vertices internal to a path p as $\text{int}(p)$. The *length* of a path p is $|p| - 1$ if the edges do not have weights and is the sum of the weights of the edges $(u, z_1), (z_1, z_2), \dots, (z_{|p|-2}, v)$ otherwise. Any path p from u to v that has the minimum length among all the paths from u to v is known as a *Shortest Path (SP)*. For any *ordered* pair of distinct vertices (u, v) , let Γ_{uv} denote the set of SPs from u to v , and let $\sigma_{uv} = |\Gamma_{uv}|$. We use $\sigma_{uv}(w)$ to denote the number of SPs from u to v that w is internal to.

The *Betweenness Centrality (BC)* $b(w)$ of a vertex $w \in V$ is [2, 29]

$$b(w) \doteq \frac{1}{n(n-1)} \sum_{\substack{(u,v) \in V \times V \\ u \neq v}} \frac{\sigma_{uv}(w)}{\sigma_{uv}} \quad (\in [0, 1]) . \quad (1)$$

There exist many variants of BC [13, 18, 26, 42, 49, 51, 52]. For ease of presentation, we focus on the one for vertices but our results can be extended to many of these variants (see Sect. 4.3).

Computing the exact BCs for all vertices is expensive [17] (see also Sect. 2). In this work, we are interested in *estimating the BCs* of all vertices, using different “approaches” (see Sect. 3.2) that offer *probabilistic guarantees* on the quality of the estimates that they compute. All approaches rely on *random sampling*: given the graph G , a user-defined *sample size* $m > 0$ and a user-defined *acceptable failure probability* $\delta \in (0, 1)$, approach A creates a *sample* $S = \{x_1, \dots, x_m\}$ by drawing m *independent* samples from an approach-specific, G -dependent, population $\mathcal{D}_{A,G}$ according to an approach-specific distribution π_A over $\mathcal{D}_{A,G}$ (for ease of notation, we will often drop G when it is clear from the context, and just use \mathcal{D}_A). A uses S to compute the estimate $\tilde{b}_{A,S}(w)$ for every vertex w (for ease of notation we will often drop S when it is clear from the

context, and just use $\tilde{b}_A(w)$). For any approach A we consider, and any $w \in V$, the estimate $\tilde{b}_{A,S}(w)$, is the *sample mean* over S of a function f_w defined over $\mathcal{D}_{A,G}$, i.e.,

$$\tilde{b}_{A,S}(w) = \frac{1}{m} \sum_{x \in S} f_w(x),$$

and it is *unbiased*, i.e., its expectation w.r.t. the choice of S is the exact BC $b(w)$. Also, for any approach A and any sample size $m > 0$, there is a function $\text{eps}_{A,m}(\cdot)$ from $\mathcal{D}_A^m \times (0, 1)$ to the non-negative reals \mathbb{R}_+ such that, with probability at least $1 - \delta$ over the runs of A with the same inputs G , m , and δ , all estimates are within $\text{eps}_{A,m}(S, \delta)$ from their exact value. In other words, it holds

$$\Pr \left(\exists w \in V \text{ s.t. } \left| b(w) - \tilde{b}_{A,S}(w) \right| > \text{eps}_{A,m}(S, \delta) \right) < \delta . \quad (2)$$

One of our goals is to compare $\text{eps}_{A_1,m}(\cdot)$, $\text{eps}_{A_2,m}(\cdot)$, ..., for different approaches A_1, A_2, \dots , as they give different characterizations of the trade-off between sample size and accuracy guarantees.

3.2 BC Estimation

We now review existing approaches for BC estimation. For each approach A , we define a domain \mathcal{D}_A , a family \mathcal{F}_A of functions from \mathcal{D}_A to $[0, 1]$, and a probability distribution π_A over \mathcal{D}_A .

The RK estimator. Riondato and Kornaropoulos [56] introduce a BC estimator that is tailored to their use of VC-dimension [63] for the analysis of the sample size sufficient to obtain a high-quality estimate of the BC of all nodes. The domain \mathcal{D}_{rk} is the set of all shortest paths between all pairs of vertices in the graph G , i.e.,

$$\mathcal{D}_{\text{rk}} \doteq \bigcup_{\substack{(u,v) \in V \times V \\ u \neq v}} \Gamma_{uv} .$$

Let $p_{uv} \in \mathcal{D}_{\text{rk}}$ be any SP from u to $v \neq u$. The distribution π_{rk} over \mathcal{D}_{rk} assigns to p_{uv} the probability mass

$$\pi_{\text{rk}}(p_{uv}) = \frac{1}{n(n-1)\sigma_{uv}} .$$

Riondato and Kornaropoulos [56] show an efficient sampling scheme to draw independent samples from \mathcal{D}_{rk} according to π_{rk} . The cost of drawing a sample is essentially that of a truncated SSSP computation from u to v . The family \mathcal{F}_{rk} of functions contains one function $f_w : \mathcal{D}_{\text{rk}} \rightarrow \{0, 1\}$ for each vertex w , defined as

$$f_w(p) \doteq \begin{cases} 1 & \text{if } w \in \text{int}(p) \\ 0 & \text{otherwise} \end{cases} .$$

The same estimator is used by Borassi and Natale [12] for a progressive sampling algorithm for BC estimation.

The ABRA estimator. The ABRA algorithm [57] uses an estimator defined over the domain

$$\mathcal{D}_{\text{ab}} \doteq \{(u, v) \in V \times V : u \neq v\},$$

i.e., over the pairs of different vertices. The distribution π_{ab} is uniform over \mathcal{D}_{ab} ; thus sampling from it is easy, per our assumptions. The family \mathcal{F}_{ab} contains a function f_w for each vertex w , defined as

$$f_w(u, v) \doteq \frac{\sigma_{uv}(w)}{\sigma_{uv}} \in [0, 1] . \quad (3)$$

Given $(u, v) \in \mathcal{D}_{\text{ab}}$, one can compute the value of $f_w(u, v)$ for each w in time proportional to performing a truncated SSSP from u to v .

The BP estimator. Jacob et al. [36], and independently Brandes and Pich [19] introduce a BC estimator that is closely related to how Brandes [17]’s algorithm computes the exact BC of all nodes. The domain \mathcal{D}_{bp} is the set V of vertices in G , and the distribution π_{bp} is uniform over this set. The family \mathcal{F}_{bp} contains one function f_w for each vertex w , defined as

$$f_w(v) \doteq \frac{1}{n-1} \sum_{z \neq v} \frac{\sigma_{vz}(w)}{\sigma_{vz}} \in [0, 1] . \quad (4)$$

The value $f_w(v)$ can be computed by performing a (full) SSSP computation from v , and then backtracking along the resulting SP DAG as in the exact BC algorithm BA by Brandes [17]. Geisberger et al. [30] present a variant of the BP estimator to ameliorate some of its issues. The theoretical guarantees of this variant and its computation are essentially the same as the BP estimator, and the same observations we make for BP can be extended to this variant.

For technical reasons, the constant zero function must belong to \mathcal{F}_{A} . This requirement is automatically satisfied when the graph contains at least one vertex with only one neighbor. When such a vertex does not exist, we just add this function to \mathcal{F}_{A} .

3.3 Bounding the Supremum Deviation

We now define Supremum Deviation, Monte-Carlo Empirical Rademacher Average, wimpy variance, and their relationship. We tailor the definitions to our settings, and refer the reader to the book by Shalev-Shwartz and Ben-David [59, Ch. 26] for in-depth discussion.

Let \mathcal{F} be a family of functions from a domain \mathcal{D} to $[0, 1]$. Let π be a distribution over \mathcal{D} , and m be a positive natural. Given a bag (sample) $\mathcal{S} = \{x_1, \dots, x_m\}$ of m samples from \mathcal{D} drawn independently according to π , the *Supremum Deviation (SD) of \mathcal{F} on \mathcal{S}* is defined as

$$\text{SD}(\mathcal{F}, \mathcal{S}) \doteq \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{\mathcal{S}}[f] - \frac{1}{m} \sum_{x \in \mathcal{S}} f(x) \right| .$$

The SD is the key concept in the study of empirical processes [53]. Sample-dependent quantities, such as the popular Empirical Rademacher average (ERA) [40], can be used to derive probabilistic upper bounds to the SD. In this work we use a Monte-Carlo estimation of the ERA, first introduced by Bartlett and Mendelson [4]. For $k \geq 1$, let $\Lambda = \{\lambda_1, \dots, \lambda_m\}$ be a bag of m i.i.d. k -dimensional Rademacher vectors, i.e., vectors whose entries are drawn independently and uniformly from $\{-1, 1\}$. The *k -Trials Monte-Carlo Empirical Rademacher Average (k -MCERA)* $\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda)$ of \mathcal{F} on \mathcal{S} using Λ is [4]

$$\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda) \doteq \frac{1}{k} \sum_{j=1}^k \left(\sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{x_i \in \mathcal{S}} \lambda_{j,i} f(x_i) \right) . \quad (5)$$

The expectation of the k -MCERA w.r.t. both \mathcal{S} and Λ is known as the *Rademacher Average* $R_m(\mathcal{F}, \pi)$ [40], and is a cornerstone of statistical learning theory. While Rademacher averages control the *expected SD*, the sharpest probabilistic bounds use the k -MCERA and depend on the *variances* of the functions in \mathcal{F} . This fact is natural, since for each $f \in \mathcal{F}$, the *variance* $\mathbb{V}[f]$ controls $\text{SD}(\{f\}, \mathcal{S})$,

asymptotically through the Central Limit Theorem, or via Bennett’s inequality [6] with finite sample guarantees. The *maximum* variance of a function $f \in \mathcal{F}$ would therefore control $\text{SD}(\mathcal{F}, \mathcal{S})$. BC values are typically very small on large graphs, thus, for f in any of the families from Sect. 3.2, it is reasonable to assume that $\mathbb{V}[f] \doteq \mathbb{E}_{\pi}[f^2] - (\mathbb{E}_{\pi}[f])^2 \approx \mathbb{E}_{\pi}[f^2]$, which is the *raw variance* of f ; henceforth all variances are raw unless otherwise noted. We define the *wimpy* variance v , and its estimator β , as

$$v \doteq \sup_{f \in \mathcal{F}} \mathbb{E}_{\pi}[f^2] \text{ and } \beta \doteq \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{x_i \in \mathcal{S}} (f(x_i))^2 . \quad (6)$$

Clearly both v and β depend on \mathcal{F} , and the families discussed in Sect. 3.2 have different wimpy variances.

The following result (proof in App. A.2) shows how to use the k -MCERA to compute an upper bound to the SD using only sample-dependent quantities.

THEOREM 3.1. *Let $\eta \in (0, 1)$, and define the quantities*

$$\begin{aligned} \gamma &\doteq \beta + \frac{2 \ln \frac{5}{\eta}}{3m} + \sqrt{\left(\frac{\ln \frac{5}{\eta}}{\sqrt{3m}} \right)^2 + \frac{2\beta \ln \frac{5}{\eta}}{m}}, \\ \rho &\doteq \hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda) + \frac{2 \ln \frac{5}{\eta}}{3km} + \sqrt{\frac{4\beta \ln \frac{5}{\eta}}{km}}, \\ r &\doteq \rho + \frac{\ln \frac{5}{\eta}}{3m} + \sqrt{\left(\frac{\ln \frac{5}{\eta}}{2\sqrt{3m}} \right)^2 + \frac{\rho \ln \frac{5}{\eta}}{m}}, \\ \varepsilon &\doteq 2r + \frac{\ln \frac{5}{\eta}}{3m} + \sqrt{\frac{2(\gamma + 4r) \ln \frac{5}{\eta}}{m}} , \end{aligned} \quad (7)$$

Then, with prob. at least $1 - \eta$ over the choice of \mathcal{S} and Λ , it holds

$$\text{SD}(\mathcal{F}, \mathcal{S}) \leq \varepsilon .$$

The quantity γ is an upper bound to v , while ρ is an upper bound to the ERA, which is the expected value of the MCERA w.r.t. Λ , and r is an upper bound to the expectation of the ERA w.r.t. \mathcal{S} . The apparent complexity of this result (in part due to the need to fit the expressions into the text column width) is due to the fact that it better characterizes the trade-off between the sample size (and other properties of the sample) and the quality of the estimates, by leveraging the *empirical wimpy variance*, i.e., the quantity β in (6). By looking at prior work under the lens of variance, one can see that the Hoeffding+union bound and the VC-dimension theorem, used respectively by Brandes and Pich [19] and Riondato and Kornaropoulos [56], assume the worst-case possible variance, while the bound used by Riondato and Upfal [57] improves to an exponential average over variances. These approaches are not sensitive to any correlation between the quantities to estimates, whereas the k -MCERA is. This advantage is one of the reasons why we choose them: this sensitivity is relevant to BC estimation, as the scores of nearby vertices are strongly correlated. Thus one should leverage these correlations when computing the approximation quality guarantee.

$$\text{It holds } \varepsilon \in 2\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda) + O\left(\frac{\ln \frac{1}{\eta}}{m} + \sqrt{\frac{(\varepsilon + R_m(\mathcal{F}, \pi)) \ln \frac{1}{\eta}}{m}} \right) .$$

Each f used for BC estimation has codomain $[0, 1]$, so $v = \mathbb{E}_\pi [f^2] \leq \mathbb{E}_\pi [f] \leq \max_{v \in V} b(v)$. Thus, ϵ is, in some sense, “output-sensitive”, as it depends on the maximum BC.

As \mathcal{F} is fixed, the quantity on the r.h.s. of (7), is a function $g_{\mathcal{F}}(\mathcal{S}, \eta)$ depending *only* on \mathcal{S} and on η . For every approach A from Sect. 3.2, we can define the function $\text{eps}_{A,m}(\cdot)$ used in (2) as

$$\text{eps}_{A,m}(\mathcal{S}, \delta) = g_{\mathcal{F}_A}(\mathcal{S}, \delta).$$

4 THE BAVARIAN FRAMEWORK

We now present BAVARIAN, our unifying algorithmic framework for BC estimation. BAVARIAN uses the k -MCERA and the variance-aware tail bound from Thm. 3.1 to compute the sample-dependent approximation quality, and can be instantiated with any of the estimators discussed in Sect. 3.2. The pseudocode is shown in Alg. 1.

Algorithm 1: BAVARIAN

```

Input: method  $A$ , graph  $G$ , sample size  $m$ , failure
      probability  $\delta$ , number of Monte Carlo trials  $k$ 
Output:  $(\tilde{B}, \epsilon)$  with the properties presented in Thm. 4.1
1  $\text{sums} \leftarrow$  map from  $V$  to vectors of size  $k+2$ 
2 foreach  $v \in V$  do  $\text{sums}[v] \leftarrow (\underbrace{-\infty, \dots, -\infty}_{k}, 0, 0)$ 
3 for  $i$  from 1 to  $m$  do
4    $s_i \leftarrow \text{drawSample}(A, G)$ 
5    $Z \leftarrow \text{getFunctionValues}(A, s_i)$ 
6    $\lambda_i \leftarrow \text{drawRademacher}(k)$ 
7   foreach  $(v, f_v(s_i)) \in Z$  do
8      $\text{sums}[v] \leftarrow \text{sums}[v] + f_v(s_i) \cdot (\lambda_{i,1}, \dots, \lambda_{i,k}, f_v(s_i), 1)$ 
9  $\tilde{B} \leftarrow \{(v, \text{sums}[v][k+2]/m), v \in V\}$ 
10  $\epsilon \leftarrow \text{getEpsilon}(\text{sums}, m, \delta)$ 
11 return  $(\tilde{B}, \epsilon)$ 
```

The input parameters to BAVARIAN are: a BC estimation method $A \in \{\text{rk}, \text{ab}, \text{bp}\}$, a graph G , a sample size $m > 0$, a failure probability $\delta \in (0, 1)$, and the number k of Monte Carlo trials for the k -MCERA (see Sect. 4.1 for a discussion of how to choose m and k). The output is a pair (\tilde{B}, ϵ) , where \tilde{B} is a set of pairs $(v, \tilde{b}_A(v))$ for each $v \in V$, where $\tilde{b}_A(v)$ is the estimate of $b(v)$ using A , and $\epsilon \in (0, 1)$ is the probabilistically-guaranteed accuracy as specified in the following theorem (proof deferred to after the description of the algorithm).

THEOREM 4.1. *With probability at least $1 - \delta$ (over the runs of the algorithm), the pair (\tilde{B}, ϵ) is such that*

$$\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \epsilon.$$

BAVARIAN first creates a map sums from V to vectors¹ of size $k+2$ (line 1 of Alg. 1), initialized to contain one key v for each $v \in V$. The vector $\text{sums}[v]$ has $k+2$ elements, with the first k initialized to $-\infty$ and the last two elements initialized to zero (line 2). We explain the role of this map below. The algorithm then enters a loop which is repeated for m iterations (lines 3–8). At the beginning of iteration i (line 4), BAVARIAN draws one element s_i from \mathcal{D}_A according to the distribution π_A . For example, if A is ab, the algorithm would draw

¹We use 1-based indexing for vectors, i.e., the first element has index 1.

a pair (u, v) of distinct vertices uniformly among all such pairs. It then uses the function `getFunctionValues` to compute, using the procedure specified by A , the set

$$Z = \{(w, f_w(s_i)) \text{ for each } w \in V \text{ s.t. } f_w(s_i) \neq 0\}.$$

For example, if A is ab, Z would contain pairs $(w, \sigma_{uv}(w)/\sigma_{uv})$ for all and only the vertices w internal to a SP from u to v , where $s_i = (u, v)$ is the population element drawn at this iteration. Together with a vector λ_i of k independent Rademacher variables (line 6), the set Z is used to update the map `sums` to maintain the following invariant: at the end of every iteration i of the loop, it must hold for each $v \in V$ that

$$\text{sums}[v][z] = \begin{cases} \sum_{j=1}^i \lambda_{j,z} f_v(s_j) & \text{for } z = 1, \dots, k \\ \sum_{j=1}^i (f_v(s_j))^2 & \text{for } z = k+1 \\ \sum_{j=1}^i f_v(s_j) & \text{for } z = k+2; \end{cases} \quad (8)$$

After m iterations of the for loop (lines 3–8), the algorithm obtains the BC estimate $\tilde{b}_A(v)$ for each vertex v as $\tilde{b}_A(v) = \text{sums}[v][k+2]/m$, and stores them in the set \tilde{B} (line 9). The probabilistic accuracy guarantee ϵ is computed by the function `getEpsilon` (line 10). This function receives as input all the necessary parameters to compute the value on the r.h.s. of (7): for each $v \in V$, it holds $\text{sums}[v][k+1] = \sum_{i=1}^k (f_v(s_i))^2 = m\beta$, and $\text{sums}[v][z] = \sum_{i=1}^k \lambda_{iz} f_v(s_i)$ for each $z = 1, \dots, k$. The output of the algorithm is then (\tilde{B}, ϵ) . We can now prove Thm. 4.1.

PROOF OF THM. 4.1. Thanks to the maintained invariant (8), each estimate $\tilde{b}_A(v)$, $v \in V$, is a sample mean (over the sample $\mathcal{S} = \{s_1, \dots, s_m\}$) of a specific function f_v , and an unbiased estimator for the exact BC $b(w) = \mathbb{E}[f_v]$. Additionally, the invariant (8) ensures that all the quantities needed to compute ϵ as in (7) are available, and thus the thesis follows from Thm. 3.1. \square

4.1 Choosing the parameters m and k

The sample size m and the number k of Monte Carlo trials are important input parameters of BAVARIAN. We now discuss how to choose m as a function of a *user-specified* desired approximation quality $\bar{\epsilon} \in (0, 1)$. In Sect. 4.1.1 discuss the choice of k as a function of m . These results are not just of theoretical interest: we use them in a *progressive-sampling* variant of BAVARIAN in Sect. 4.2.

Minimum sample size. We start by deriving an upper bound to the minimum sample size m needed for the approximation quality guarantee ϵ returned by the algorithm to be *no-greater* than a user-desired upper bound $\bar{\epsilon}$. From Thm. 3.1 it is evident that ϵ decreases as $\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \lambda)$ and β from (6) decrease, and as m and k increase. Consider a procedure `zeroEpsilon`(m, k, η) that, given m, k , and η , computes ϵ as in Thm. 3.1 with the assumption that $\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \lambda) = 0$ (which is the minimum possible because of the requirement that the zero constant function belongs to \mathcal{F}), and $\beta = 0$, and let

$$m_*(k, \eta, \bar{\epsilon}) = \min\{m > 0 : \text{zeroEpsilon}(m, k, \eta) \leq \bar{\epsilon}\}. \quad (9)$$

FACT 1. *BAVARIAN would never return a pair (\tilde{B}, ϵ) with $\epsilon \leq \bar{\epsilon}$ if the passed sample size m is smaller than $m_*(k, \eta, \bar{\epsilon})$ (and the other input parameters are k and $\delta = \eta$).*

Computing $m_*(k, \eta, \bar{\epsilon})$ is easily done with a binary search on m using `zeroEpsilon`.

For $m \geq m_*(k, \eta, \bar{\varepsilon})$, the algorithm *may* return $\varepsilon < \bar{\varepsilon}$. We leverage this property in Sect. 4.2, where we use (9) to derive a sample schedule for a progressive-sampling variant of BAVARIAN.

Sufficient sample size. We now present a result connecting structural properties of the graph G to a sample size *sufficient* for the algorithm to return an approximation quality bound ε not larger than the user-desired $\bar{\varepsilon}$. In the general case, the sample size we show depends on the *vertex-diameter* $\text{vd}(G)$ of the graph G [56, Def. 1], i.e., the *maximum number of vertices on a SP in G* . On unweighted graphs, the vertex diameter equals the diameter plus one, as $\text{vd}(G)$ counts *vertices*, while the diameter counts *edges*. On weighted graphs, the two quantities are not necessarily related. All that is really needed is an *upper bound* to the vertex-diameter, which is easy to compute in either case (see below). Riondato and Kornaropoulos [56, Coroll. 1] showed that the *Vapnik-Chervonenkis (VC-) dimension* [63] of the task of approximating BC using the rk estimator is upper bounded by $\lfloor \log_2(\text{vd}(G) - 2) \rfloor + 1$. They then use this result to derive a sufficient sample size for their BC approximation algorithm. They also show that, when G is undirected and there is no more than one SP between any pair of vertices in G , as is sometimes enforced on road networks, the VC-dimension collapses to 3, independently from the vertex-diameter of the graph [56, Lemma 2]. These results are limited to the rk estimator and do not make use of the Rademacher averages, which could lead to a much smaller sufficient sample size. We *extend and adapt* these results so that they can also be used to derive sufficient sample sizes for the bp and ab estimators and when using Rademacher averages, as shown in the following theorem.

THEOREM 4.2. *Let $d \geq \text{vd}(G)$ or $d = 3$ if G is undirected and there is at most one SP between any pair of vertices in G . Let*

$$m^*(\bar{\varepsilon}, \eta) \doteq \frac{1}{\bar{\varepsilon}^2} \left(4Cd + 4\sqrt{\frac{Cd \ln \frac{2}{\eta}}{2} + \frac{\ln \frac{2}{\eta}}{2}} \right), \quad (10)$$

where C is a universal constant [33, 41, 60]. When run with inputs $m \geq m^*(\bar{\varepsilon}, \eta)$, and $\delta = \eta$, BAVARIAN could return $(\tilde{B}, \bar{\varepsilon})$ rather than using the computed ε , and Thm. 4.1 would still hold.

The proof of this theorem (in App. A.2) relies on a novel family of functions that we call *SP-enforcing family* (Def. A.7), which we deem of independent interest, and on a result (Lemma A.5) to upper bound the ERA, i.e., the expectation w.r.t. λ of the k -MCERA, with the VC-dimension, used to adapt the result on the VC-dimension of BC-estimation with the rk estimator [56, Coroll. 1, Lemma 2] to Rademacher averages. Theorem 4.2 is *not only of theoretical interest*. Rather, it has a double *practical impact*: (1) a better characterization of the relationship between sample size and accuracy guarantee enables the user to better select what sample size to use; (2) pushing down the sample size sufficient to obtain a desired approximation guarantee directly benefits progressive sampling algorithms (see Sect. 4.2), which can now deterministically terminate earlier.

It is evident from the statement of Thm. 4.2 that an *upper bound* to the vertex-diameter of G is sufficient to compute the sufficient sample size. In unweighted graphs, it is possible to obtain a 2-approximation of the diameter (thus of the vertex-diameter) with one SSSP computation by summing the lengths of the two longest

SPs. Bergamini and Meyerhenke [7, Sect. 3] show how to get a good upper bound to $\text{vd}(G)$ on weighted graphs.

4.1.1 Choosing the number of Monte-Carlo trials. The beauty of the Monte-Carlo bound ρ in Thm. 3.1 is that the uncertainty error term added to $\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda)$ is *asymptotically negligible* for any choice of k (i.e., k is absent from the asymptotic form). A value $k \geq 8$ would match the $8r$ term inside the square root of ε , after which the Monte-Carlo estimation terms are negligible compared to other error terms. Constant-factor improvements are still possible for $k > 8$, though due to rapidly diminishing returns in k , additional computational resources are better spent on computing more shortest paths (i.e., taking more samples) than increasing k unboundedly.

From a computational perspective, k can be selected so the cost of computing SPs amortizes the cost of running k Monte-Carlo trials. Computing the MCERA requires $O(mk|V|)$ time for BC (often much less for the rk and ab estimators, see also Sect. 5). Assuming BFS with time complexity $O(|E|)$, it makes sense to have $k \in O(|E|/|V|)$ trials, and with Dijkstra's algorithm (time complexity $O(|E| + |V| \log |V|)$), to have $k \in O(|E|/|V| + \log |V|)$.

4.2 Progressive sampling algorithm

The algorithmic framework discussed so far uses *static* or *one-shot* sampling: it draws a *single* sample \mathcal{S} of user-specified size m from the appropriate population, and uses information from the sample to compute the quality guarantee ε . In Sect. 4.1 we discussed how the user can make an informed decision on m based on their desired approximation quality $\bar{\varepsilon}$. The sample sizes m_* and m^* from (9) and (10) give a range $[m_*, m^*]$ such that BAVARIAN may return a value ε smaller than $\bar{\varepsilon}$ when the given sample size m is within this range. The likelihood that ε is smaller than $\bar{\varepsilon}$ clearly increases as m gets closer to m^* , and the algorithm is *guaranteed* to return $\varepsilon < \bar{\varepsilon}$ when $m \geq m^*$. We now introduce a more flexible variant of BAVARIAN based on *progressive sampling*: rather than using a single fixed size sample, our algorithm BAVARIAN-P uses a *sample-schedule* $\Delta = (m_i)_{i=0}^{T-1}$ to *iteratively grow* the sample \mathcal{S} from an initial size m_0 to larger sizes m_1, \dots, m_{T-1} (where T is the maximum number of iterations, discussed below). The algorithm stops at the earliest iteration i such that the approximation quality ε obtained from the sample \mathcal{S} at iteration i is no larger than the user-desired quality $\bar{\varepsilon}$. Thus, $\varepsilon \leq \bar{\varepsilon}$ is the *stopping condition* of our algorithm. The computation of ε must handle the fact that we are essentially analyzing multiple samples, as discussed in the proof of Thm. 4.3, which states the properties of BAVARIAN-P. We first present the algorithm (pseudocode in Alg. 2) with the sample schedule Δ and the desired quality $\bar{\varepsilon}$ passed as an input parameters, and then discuss appropriate choices for Δ . The only requirement on Δ is

$$m_{T-1} \geq m^*(\bar{\varepsilon}, \delta/T) \quad (11)$$

for technical reasons explained in the proof of Thm. 4.3.

BAVARIAN-P first initializes the data structure sums exactly as in the static-sampling case (see Alg. 1 and its description), and then, starting at $i = 0$, enters a loop (lines 5–10) that continues until the stopping condition is satisfied (more details in the following). At each iteration of the loop, the function `drawSamplesAndUpdateSums` is called (line 6). It performs the operations on lines 3 to 8 of Alg. 1: it samples from the population (for $m_i - m_{i-1}$ times), computes the

Algorithm 2: BAVARIAN-P

Input: method A, graph G , sample schedule $\Delta = (m_i)_{i=0}^{T-1}$, failure prob. δ , no. of MC-trials k , desired quality $\bar{\varepsilon}$

Output: (\tilde{B}, ε) with the properties presented in Thm. 4.3

- 1 sums \leftarrow map from V to vectors of size $k + 2$
- 2 **foreach** $v \in V$ **do** sums[v] \leftarrow $\underbrace{(-\infty, \dots, -\infty)}_k, 0, 0)$
- 3 $m_{-1} \leftarrow 0$
- 4 $i \leftarrow 0$
- 5 **do**
- 6 | drawSamplesAndUpdateSums(sums, $m_i - m_{i-1}$, A, G)
- 7 | **if** $i < T - 1$ **then** $\varepsilon \leftarrow \text{getEpsilon}(\text{sums}, m_i, \delta/T)$
- 8 | **else** $\varepsilon \leftarrow \bar{\varepsilon}$
- 9 | $i \leftarrow i + 1$
- 10 **while** $\varepsilon > \bar{\varepsilon}$
- 11 $\tilde{B} \leftarrow \{(v, \text{sums}[v][k+2]/m_i), v \in V\}$
- 12 **return** (\tilde{B}, ε)

function values, draws a Rademacher vector, and appropriately updates sums while maintaining the invariant from (8). Then (line 7), the algorithm computes the approximation quality guarantee ε that can be obtained from the sample of size m_i seen so far, using δ/T as the failure probability (as discussed for Alg. 1). BAVARIAN-P stops iterating when the computed ε is not larger than the desired $\bar{\varepsilon}$, which is guaranteed to be the case when i equals $T - 1$, due to (11). The approximation \tilde{B} is computed (line 11) and output together with ε . The proof of the following result is in App. A.2.

THEOREM 4.3. *With probability at least $1 - \delta$ (over the runs of the algorithm), the pair (\tilde{B}, ε) returned by BAVARIAN-P is such that*

$$\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \varepsilon \leq \bar{\varepsilon} .$$

Choosing the sample schedule. We now discuss reasonable choices for the sample schedule $\Delta = (m_i)_{i=0}^{T-1}$. Assume, for ease of presentation, that the number of iterations T is a fixed parameter chosen by the user. We later remove this assumption. The correctness of the algorithm requires that the last sample size m_{T-1} satisfies (11) and, given what we discussed in Sect. 4.1, it makes no sense to use a larger sample size than the quantity in the r.h.s. of (11). Thus, we assume that m_{T-1} equals this quantity. It follows from Fact 1 that it also does not make sense to have $m_0 \leq m_*(k, \delta/T, \bar{\varepsilon})$, where m_* is defined in (9). It is not necessarily appropriate to use *exactly* this sample size as the first one, for the reason argued in Sect. 4.1, but this quantity acts as a lower bound. Once m_0 and m_{T-1} have been fixed, we need to choose $T - 2$ intermediate sample sizes between these two extremes. A reasonable approach is to follow a *geometrically-increasing* sample schedule, which has been shown in practice and in theory to be close to optimal [37]. Thus, we have

$$m_i = m_0 \left(\frac{m_{T-1}}{m_0} \right)^{\frac{i}{T-1}}, \text{ for } i = 1, \dots, T - 2 .$$

The scaling factor m_{T-1}/m_0 is a result of fixing the number of iterations in advance. It is more convenient to let the user specify a desired scaling factor $\theta > 1$. An unbounded binary search then finds the minimum number of iterations T^* such that there are

exactly $T^* - 2$ sample sizes $m_i = m_0 \theta^i$ for $i = 1, \dots, T^* - 2$ between the smallest sample size (computed using (9) with $\eta = \delta/T^*$) and the largest m_{T^*-1} (from (11), potentially smaller than $m_0 \theta^{T^*-1}$).

4.3 Extensions

The unifying approach taken by BAVARIAN and BAVARIAN-P allows them to be adapted to work with many variants of BC and with other centrality measures. Due to space limitations, we only remark that the extensions to other BC variants would follow the same approach discussed in [56, Sect. 6]. de Lima et al. [25] showed how to estimate the percolation centrality through sampling by adapting the work of Riondato and Kornaropoulos [56] and Riondato and Upfal [57]. One can similarly adapt BAVARIAN to estimate percolation centrality. When dealing with fully-dynamic graphs with arbitrary insertions and/or deletions of edges and/or vertices, BAVARIAN can be used as a drop-in replacement for ABRA as discussed in [57, Sect. 5].

An interesting direction for future work is understanding the requirements for centrality measures to be amenable to approximation using the MCERAs and/or similar tools. The extension is less straightforward than it may seem at first, due to the need for efficient sampling schemes from the appropriate populations.

5 EXPERIMENTAL EVALUATION

We now present the results of our experimental evaluation of BAVARIAN and of the comparison between the different estimators presented in Sect. 3.2. We do not discuss the performance of BAVARIAN-P separately, as it is tightly coupled with that of BAVARIAN. We report here only a subset of the results. Results for all datasets and parameters are in App. A.3.

Goals. A first goal is to evaluate the behavior of the quality guarantee ε output by BAVARIAN as the input parameters m and k change, and to compare it with that of the state of the art for each estimator. This comparison allows us to assess the power of the MCERAs and of the variance-aware tail bound on the SD (Thm. 3.1). We also evaluate the trade-off between ε and runtime for the different BC estimators, to understand which one is more “efficient” in terms of the accuracy guarantee improvement per unit of time. BAVARIAN *does not in any way change the empirical properties of the estimators* (e.g., the similarity between the ranking of the vertices by estimated BC and the ranking by exact BC), so we do not evaluate them. For in-depth evaluations of such very important aspects, which are not impacted by BAVARIAN, see the original works [19, 56, 57] and the experimental comparisons by Matta et al. [48] and AlGhamdi et al. [1].

Table 1: Dataset characteristics

Graph	Vertices	Edges
ca-AstroPh	17,903	197,031
cit-HepPh	34,546	421,578
email-Enron	36,682	183,831
p2p-Gnutella31	62,586	147,892
soc-Epinions1	75,879	508,837
com-dblp	317,080	1,049,866
com-youtube	1,134,890	2,987,624

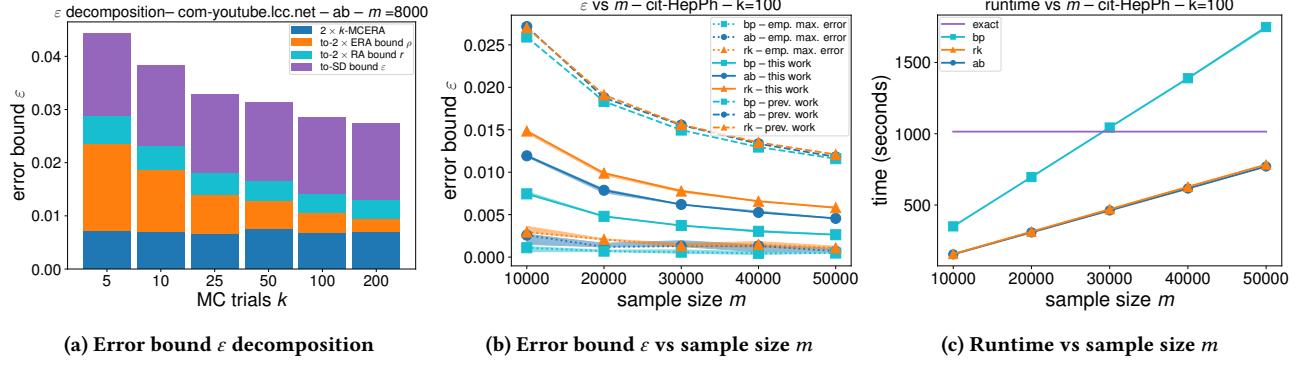


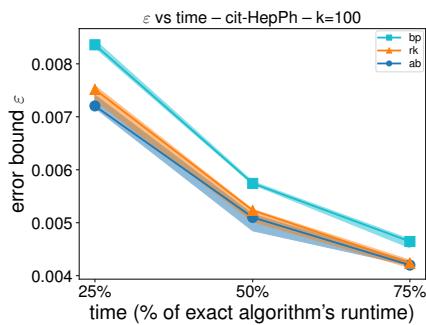
Figure 1: Experimental results. See description in the text.

Implementation, environment, and datasets. We implemented BAVARIAN in C++20 as an extension of NetworkKit [61], and compared it with the implementation of Brandes’s algorithm [17] included in the same suite. The code for reproducing all experiments and figures is in the same archive where you found this extended version. We run them on a FreeBSD 11 Amazon AWS instance with 8 AMD EPYC CPUs and 32GB of RAM. We used different combinations of values for m , k , δ , and on different graphs (see below), doing five runs per combination. The difference across runs is minimal, which is not surprising as most of our results are not *absolutely* tight, while being *tighter* than previous works. We report the results for the median, the maximum and the minimum over the runs (represented as shaded areas around the curve for the median). We only show results for $\delta = 0.1$. This parameter has minimal impact on the performances of our algorithm, and results for other values of δ are qualitatively similar.

We used datasets from the SNAP repository [44]. Their salient characteristics are reported in Table 1.

5.1 Results

We start by remarking that in all the hundreds of runs we performed, the supremum deviation (i.e., the maximum error in our estimation) was always less than the error ε returned by BAVARIAN, and often significantly so. This fact is not surprising, as Thm. 3.1 is not necessarily tight. Hence, BAVARIAN is even more accurate than the theory guarantees, which leaves space for future improvements.

Figure 2: Error bound ε vs time.

Impact of k . We evaluate the impact of k on ε , as the impact of k on the runtime is negligible, because BAVARIAN’s runtime is heavily dominated by the SSSP computations. We run BAVARIAN for different values of k and split ε in its components from Thm. 3.1. In Fig. 1a we show the decomposition, and observe the behavior as function of k , on the com-youtube graph for the ab estimator, with $m = 8000$ (results for other values of the parameters, other datasets, and other estimators are qualitatively similar). The x -axis of the figure has essentially a logarithmic scale. We can see that ε rapidly decreases as a function of k , but the returns are diminishing. The decrease of ε is readily explained by looking at its components: the dark blue (bottom) part of the bars is *twice* the k -MCERA (i.e., $2\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda)$), the orange (2^{nd} from bottom) portion is the difference between 2ρ and $2\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda)$, the cyan (3^{rd} from bottom) is the difference between $2r$ and 2ρ , and the purple (top portion) is the difference between ε and $2r$. The only part that changes significantly with k is the orange one, which corresponds to the probabilistic tail bound from the k -MCERA to the ERA, its expectation w.r.t. Λ . As the k -MCERA is tightly concentrated around the ERA, it is not surprising that with even a relatively small number of MC trials, this tail bound becomes negligible w.r.t. the other terms.

Impact of m . Figures 1b and 1c show the impact of the sample size m on the error bound ε and on the runtime, for cit-HepPh with $k = 100$ (results for other datasets/parameters are qualitatively similar). As expected, the error bound decreases approximately as $1/\sqrt{m}$ and the runtime increases linearly with m . We can also observe that the use of the k -MCERA and of Thm. 3.1 allows BAVARIAN to obtain much smaller error bounds than previous works for all estimators, and much closer to the empirical maximum error measured on the sample. There is still room for improvement in tightening the error bounds, but the leap forward is significant. Looking only at Fig. 1b one may be tempted to conclude that the bp estimator should be preferred over the others, because it results in a smaller ε . This conclusion would not be appropriate: the estimators sample from different populations and do different work per sample, so their performance at equal sample sizes is not comparable. By also looking at Fig. 1c, it is evident that BAVARIAN with bp takes much longer at each sample size than with other estimators (which take approximately the same time because they do similar work per sample), and becomes slower than the exact algorithm earlier. Thus,

per sample, bp is more “efficient” than the other estimators, but it is *less efficient per unit of time*. Fig. 2 shows how, when run for the same amounts of time, bp gives a much worse (i.e., larger) ε than the other estimators. We conclude from this comparison that the ab estimator gives the best “bang for the buck”, i.e., the best ε per unit of time, and should likely be preferred over the others.

6 CONCLUSIONS

We present BAVARIAN, a suite of algorithms for approximating the BC of all vertices in a graph with stringent probabilistic quality guarantees. BAVARIAN achieves a better trade-off between sample size and approximation quality thanks to the use of the k -MCERAs and their variance-aware bounds. The unifying framework offered by BAVARIAN allows for a fair comparison of different sampling-based estimators for BC, and for us to generalize sample-complexity results that held only for one estimator. The results of our experimental evaluation of BAVARIAN show it guarantees much better approximations at smaller sample sizes, i.e., faster, than the previous state of the art. A few MC-trials (i.e., low k) are sufficient to get a low error bound. The ab estimator stands out as more efficient than the others in terms of “unit of guarantee” per unit of time. In future work, we plan to generalize the application of the MCERAs to other centrality measures, and to design more efficient progressive sampling algorithms, leveraging martingales approaches.

ACKNOWLEDGMENTS

Part of this work is supported by the National Science Foundation grant 2006765 and by the DARPA/ARFL grant FA8750.

REFERENCES

- [1] Ziyad AlGhamdi, Fuad Jamour, Spiros Skiadopoulos, and Panos Kalnis. 2017. A Benchmark for Betweenness Centrality Approximation Algorithms on Large Graphs. In *Proceedings of the 29th International Conference on Scientific and Statistical Database Management (SSDBM '17)*. ACM.
- [2] Jac. M. Anthonisse. 1971. *The rush in a directed graph*. Technical Report BN 9/71. Stichting Mathematisch Centrum, Amsterdam, Netherlands.
- [3] David A. Bader, Shiva Kintali, Kamesh Madduri, and Milena Mihail. 2007. Approximating Betweenness Centrality. In *Algorithms and Models for the Web-Graph*, Anthony Bonato and FanR.K. Chung (Eds.). Lecture Notes in Computer Science, Vol. 4863. Springer Berlin Heidelberg, 124–137.
- [4] Peter L. Bartlett and Shahar Mendelson. 2002. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research* 3, Nov (2002), 463–482.
- [5] Alex Bavelas. 1950. Communication patterns in task-oriented groups. *The Journal of the Acoustical Society of America* 22, 6 (1950), 725–730.
- [6] George Bennett. 1962. Probability inequalities for the sum of independent random variables. *J. Amer. Statist. Assoc.* 57, 297 (1962), 33–45.
- [7] Elisabetta Bergamini and Henning Meyerhenke. 2015. Fully-dynamic Approximation of Betweenness Centrality. In *Proceedings of the 23rd European Symposium on Algorithms (ESA '15)*, 155–166.
- [8] Elisabetta Bergamini and Henning Meyerhenke. 2016. Approximating Betweenness Centrality in Fully-dynamic Networks. *Internet Mathematics* 12, 5 (2016), 281–314.
- [9] Elisabetta Bergamini, Henning Meyerhenke, and Christian L. Staudt. 2015. Approximating Betweenness Centrality in Large Evolving Networks. In *17th Workshop on Algorithm Engineering and Experiments (ALENEX '15)*. SIAM, 133–146.
- [10] Paolo Boldi and Sebastiano Vigna. 2014. Axioms for centrality. *Internet Mathematics* 10, 3-4 (2014), 222–262.
- [11] Francesco Bonchi, Gianmarco De Francisci Morales, and Matteo Riondato. 2016. Centrality measures on big graphs: Exact, approximated, and distributed algorithms. In *Proceedings of the 25th international conference companion on world wide web*, 1017–1020.
- [12] Michele Borassi and Emanuele Natale. 2019. KADABRA is an ADaptive Algorithm for Betweenness via Random Approximation. *Journal of Experimental Algorithms* 24, 1 (2019), 1–35.
- [13] Stephen P. Borgatti and Martin G. Everett. 2006. A Graph-theoretic perspective on centrality. *Soc. Netw.* 28, 4 (2006), 466–484.
- [14] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. 2000. A sharp concentration inequality with application. *Random Structures & Algorithms* 16, 3 (2000), 277–292.
- [15] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. 2013. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- [16] Olivier Bousquet. 2002. A Bennett concentration inequality and its application to suprema of empirical processes. *Comptes Rendus Mathématique* 334, 6 (2002), 495–500.
- [17] Ulrik Brandes. 2001. A faster algorithm for betweenness centrality. *Journal of Mathematical Sociology* 25, 2 (2001), 163–177.
- [18] Ulrik Brandes. 2008. On variants of shortest-path betweenness centrality and their generic computation. *Social Networks* 30, 2 (2008), 136–145.
- [19] Ulrik Brandes and Christian Pich. 2007. Centrality estimation in large networks. *International Journal of Bifurcation and Chaos* 17, 7 (2007), 2303–2318.
- [20] Sergio Cabello, Erin W. Chambers, and Jeff Erickson. 2013. Multiple-source shortest paths in embedded graphs. *SIAM J. Comput.* 42, 4 (2013), 1542–1571.
- [21] Mostafa Haghif Chehreghani, Albert Bifet, and Talel Abdessalem. 2018. Efficient Exact and Approximate Algorithms for Computing Betweenness Centrality in Directed Graphs. In *Advances in Knowledge Discovery and Data Mining*, Dinh Phung, Vincent S. Tseng, Geoffrey I. Webb, Bao Ho, Mohades Ganji, and Lida Rashidi (Eds.). Springer International Publishing, Cham, 752–764.
- [22] Flavio Chierichetti, Anirban Dasgupta, Ravi Kumar, Silvio Lattanzi, and Tamás Sarlós. 2016. On sampling nodes in a network. In *Proceedings of the 25th International Conference on World Wide Web*, 471–481.
- [23] Flavio Chierichetti and Shahrad Haddadan. 2018. On the Complexity of Sampling Vertices Uniformly from a Graph. In *45th International Colloquium on Automata, Languages, and Programming (ICALP 2018)*.
- [24] Cyrus Cousins and Matteo Riondato. 2020. Sharp uniform convergence bounds through empirical centralization. In *Advances in Neural Information Processing Systems*.
- [25] Alane M. de Lima, Murilo V. G. da Silva, and André L. Vignati. 2020. Estimating the Percolation Centrality of Large Networks through Pseudo-dimension Theory. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM.
- [26] Shlomi Dolev, Yuval Elovici, and Rami Puzis. 2010. Routing betweenness centrality. *J. ACM* 57, 4, Article 25 (May 2010), 27 pages.
- [27] Dóra Erdős, Vatche Ishakian, Azer Bestavros, and Evinaria Terzi. 2015. A Divide-and-Conquer Algorithm for Betweenness Centrality. In *SIAM International Conference on Data Mining (SDM '15)*. SIAM, 433–441.
- [28] Changjun Fan, Li Zeng, Yuhui Ding, Muhan Chen, Yizhou Sun, and Zhong Liu. 2019. Learning to Identify High Betweenness Centrality Nodes from Scratch. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. ACM. <https://doi.org/10.1145/3357384.3357979>
- [29] Linton C. Freeman. 1977. A set of measures of centrality based on betweenness. *Sociometry* 40 (1977), 35–41.
- [30] Robert Geisberger, Peter Sanders, and Dominik Schultes. 2008. Better Approximation of Betweenness Centrality. In *10th Workshop on Algorithm Engineering and Experiments (ALENEX '08)*. SIAM, 90–100.
- [31] Jay Ghurye and Mihai Pop. 2016. Better Identification of Repeats in Metagenomic Scaffolding. In *WABI 2016: Algorithms in Bioinformatics*. Springer, 174–184.
- [32] Oded Green, Robert McCall, and David A. Bader. 2012. A Fast Algorithm for Streaming Betweenness Centrality. In *2012 International Conference on Privacy, Security, Risk and Trust (PASSAT '12)*. IEEE, 11–20.
- [33] David Haussler. 1995. Sphere packing numbers for subsets of the Boolean n-cube with bounded Vapnik-Chervonenkis dimension. *Journal of Combinatorial Theory, Series A* 69, 2 (1995), 217–232.
- [34] Takanori Hayashi, Takuya Akiba, and Yuichi Yoshida. 2015. Fully Dynamic Betweenness Centrality Maintenance on Massive Networks. *Proceedings of the VLDB Endowment* 9, 2 (2015), 48–59.
- [35] Wassily Hoeffding. 1963. Probability Inequalities for Sums of Bounded Random Variables. *J. American Statistical Assoc.* 58, 301 (1963), 13–30.
- [36] Riko Jacob, Dirk Koschützki, Katharina Anna Lehmann, Leon Peeters, and Dagmar Tenfelde-Podehl. 2005. Algorithms for Centrality Indices. In *Network Analysis*, Ulrik Brandes and Thomas Erlebach (Eds.). Lecture Notes in Computer Science, Vol. 3418. Springer Berlin Heidelberg, 62–82.
- [37] George H. John and Pat Langley. 1996. Static Versus Dynamic Sampling for Data Mining. In *Proc. 2nd Int. Conf. Knowl. Disc. Data Mining (KDD '96)*. The AAAI Press, Menlo Park, CA, USA, 367–370.
- [38] Miray Kas, Matthew Wachs, Kathleen M. Carley, and L. Richard Carley. 2013. Incremental Algorithm for Updating Betweenness Centrality in Dynamically Growing Networks. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM '13)*. IEEE/ACM, 33–40.
- [39] Liran Katzir, Edo Liberty, Oren Somekh, and Ioana A Cosma. 2014. Estimating sizes of social networks via biased sampling. *Internet Mathematics* 10, 3-4 (2014), 335–359.

- [40] Vladimir Koltchinskii. 2001. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory* 47, 5 (July 2001), 1902–1914.
- [41] Aryeh Kontorovich. 2016. Agnostic PAC lower bound. <https://www.cs.bgu.ac.il/~asm162/wiki/files/agnostic-pac-lb.pdf>
- [42] Nicolas Kourtellis, Tharaka Alahakoon, Ramanuja Simha, Adriana Iamnitchi, and Rahul Tripathi. 2012. Identifying high betweenness centrality nodes in large social networks. *Social Network Analysis and Mining* 3, 4 (2012), 899–914.
- [43] Nicolas Kourtellis, Gianmarco De Francisci Morales, and Francesco Bonchi. 2015. Scalable Online Betweenness Centrality in Evolving Graphs. *IEEE Transactions on Knowledge and Data Engineering* 27, 9 (2015), 2494–2506.
- [44] Jure Leskovec and Andrej Krevl. 2014. SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>.
- [45] Yixia Li, Shudong Li, Yanshan Chen, Peiyuan He, Xiaobo Wu, and Weihong Han. 2019. Electric Power Grid Invulnerability Under Intentional Edge-Based Attacks. In *Dependsys 2019: Dependability in Sensor, Cloud, and Big Data Systems and Applications*. Springer Singapore, 454–461.
- [46] Yeon-sup Lim, Daniel S. Menasche, Bruno Ribeiro, Don Towsley, and Prithwish Basu. 2011. Online estimating the k central nodes of a network. In *IEEE Network Science Workshop (NSW'11)*, 118–122.
- [47] Arun S. Maiya and Tanya Y. Berger-Wolf. 2010. Online Sampling of High Centrality Individuals in Social Networks. In *Advances in Knowl. Disc. Data Mining*. Springer Berlin Heidelberg, 91–98.
- [48] John Matta, Gunes Ercal, and Koushik Sinha. 2019. Comparing the speed and accuracy of approaches to betweenness centrality approximation. *Computational Social Networks* 6, 1 (2019), 2.
- [49] Adam McLaughlin and David A. Bader. 2014. Scalable and High Performance Betweenness Centrality on the GPU. *SC14: International Conference for High Performance Computing, Networking, Storage and Analysis* (Nov 2014).
- [50] Mark E. J. Newman and Michelle Girvan. 2004. Finding and evaluating community structure in networks. *Phys. Rev. E* 69 (Feb. 2004). Issue 2.
- [51] Tore Opsahl, Filip Agneessens, and John Skvoretz. 2010. Node centrality in weighted networks: Generalizing degree and shortest paths. *Soc. Netw.* 32, 3 (2010), 245–251.
- [52] Jürgen Pfeffer and Kathleen M. Carley. 2012. k-Centrailities: local approximations of global measures based on shortest paths. In *Proc. 21st Int. Conf. Companion on World Wide Web (WWW '12 Companion)*. ACM, New York, NY, USA, 1043–1050.
- [53] David Pollard. 1984. *Convergence of stochastic processes*. Springer-Verlag.
- [54] Matteo Pontecorvi and Vijaya Ramachandran. 2015. Fully Dynamic Betweenness Centrality. In *Proceedings of the 26th International Symposium on Algorithms and Computation (ISAAC '15)*, 331–342.
- [55] Dimitrios Proutzlos and Keshav Pingali. 2013. Betweenness centrality: algorithms and implementations. In *Proc. 18th ACM SIGPLAN Symp. Principles and Practice of Parallel Programming* (Shenzhen, China) (PPoPP '13). ACM, New York, NY, USA, 35–46.
- [56] Matteo Riondato and Evangelos M. Kornaropoulos. 2016. Fast approximation of betweenness centrality through sampling. *Data Mining and Knowledge Discovery* 30, 2 (2016), 438–475.
- [57] Matteo Riondato and Eli Upfal. 2018. ABRA: Approximating Betweenness Centrality in Static and Dynamic Graphs with Rademacher Averages. *ACM Trans. Knowl. Disc. from Data* 12, 5 (2018), 61.
- [58] Ahmet Erdem Sarıyüce, Kamer Kaya, Erik Saule, and Ümit V. Çatalyürek. 2017. Graph Manipulations for Fast Centrality Computation. *ACM Transactions on Knowledge Discovery from Data* 11, 3 (2017), 1–25.
- [59] Shai Shalev-Shwartz and Shai Ben-David. 2014. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press.
- [60] Nathan Srebro and Karthik Sridharan. 2010. Note on refined Dudley integral covering number bound. (2010). <http://www.cs.cornell.edu/~sridharan/dudley.pdf>.
- [61] Christian L. Staudt, Aleksejs Sazonovs, and Henning Meyerhenke. 2016. NetworkKit: An Interactive Tool Suite for High-Performance Network Analysis. *Network Science* 4, 4 (2016).
- [62] Vladimir N. Vapnik. 1998. *Statistical learning theory*. Wiley.
- [63] Vladimir N. Vapnik and Alexey J. Chervonenkis. 1971. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability and its Applications* 16, 2 (1971), 264–280.
- [64] Yuichi Yoshida. 2014. Almost Linear-time Algorithms for Adaptive Betweenness Centrality Using Hypergraph Sketches. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New York, New York, USA) (KDD '14). ACM, New York, NY, USA, 1416–1425.

A ADDITIONAL MATERIALS FOR BAVARIAN

A.1 Reproducibility and Additional Results

All the code and datasets can be found in the same archive where you found this file. It includes a `README.txt` file with instructions,

and a script to run all our experiments and to generate all the figures. Additional experimental results can be found in App. A.3 below.

A.2 Missing Proofs

We now seek to prove Thm. 3.1. In service of this result, we require a few intermediaries.

Definition A.1. A function $Z \in \mathcal{D}^m \rightarrow \mathbb{R}$ is *self-bounding* with scale $\gamma \geq 0$, for some $\gamma \geq 0$ if for each $j = 1, \dots, m$, there exists a function $Z_j \in \mathcal{D}^m \rightarrow \mathbb{R}$ such that, for any $\mathcal{S} \in \mathcal{D}^m$ it holds that

- (1) $Z_j(\mathcal{S})$ does not depend on the j -th component \mathcal{S}_j of \mathcal{S} ; and
- (2) it holds $Z_j(\mathcal{S}) \leq Z(\mathcal{S}) \leq Z_j(\mathcal{S}) + \gamma$;

Additionally, the functions Z_j , $j = 1, \dots, m$, must be such that, for any $\mathcal{S} \in \mathcal{D}^m$, it holds $\sum_{j=1}^m (Z(\mathcal{S}) - Z_j(\mathcal{S})) \leq Z(\mathcal{S})$.

THEOREM A.2. Let Z be a function from \mathcal{D}^m to \mathbb{R} that is self-bounding with scale γ . Let $\delta \in (0, 1)$, and \mathcal{S} be a collection of m i.i.d. samples from \mathcal{D} . With prob. $\geq 1 - \delta$ over the choice of \mathcal{S} , it holds

$$\mathbb{E}_{\mathcal{S}} [Z(\mathcal{S})] \leq Z(\mathcal{S}) + \frac{2\gamma}{3} \ln \frac{1}{\delta} + \sqrt{\left(\frac{\gamma}{\sqrt{3}} \ln \frac{1}{\delta} \right)^2 + 2\gamma Z(\mathcal{S}) \ln \frac{1}{\delta}} . \quad (12)$$

This result is derived using the *sub-gamma* form in [15, Sect. 2.1] of the stronger *sub-Poisson* self-bounding function inequality [14, Thm. 1]. This form appeared as [24, Thm. 6].

THEOREM A.3. With probability at least $1 - \eta$ over the choice of \mathcal{S} , it holds that

$$\nu \leq \beta + \frac{2 \ln \frac{5}{\eta}}{3m} + \sqrt{\left(\frac{\ln \frac{5}{\eta}}{\sqrt{3m}} \right)^2 + \frac{2\beta \ln \frac{5}{\eta}}{m}} .$$

PROOF. We show this result by showing that an upper-bound on the raw wimpy variance is self-bounding with scale $\frac{1}{m}$. First note that by Jensen's inequality,

$$m\nu = \sup_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{S}} \left[\sum_{i=1}^m (f(\mathcal{S}_i))^2 \right] \leq \mathbb{E}_{\mathcal{S}} \left[\sup_{f \in \mathcal{F}} (f(\mathcal{S}_i))^2 \right] = m\mathbb{E}_{\mathcal{S}} [\beta] .$$

The remainder of this proof thus seeks to bound the r.h.s., as this yields a bound on ν , as desired.

Suppose $\mathcal{S} \in \mathcal{D}^m$, let $\mathcal{S}_{\setminus j}$ denote the vector \mathcal{S} , missing the j th component, note that

$$Z(\mathcal{S}) = \sup_{f \in \mathcal{F}} \sum_{i=1}^m f^2(\mathcal{S}_i) ,$$

and take

$$Z_j(\mathcal{S}) \doteq \sup_{f \in \mathcal{F}} \sum_{i=1, i \neq j}^m f^2(\mathcal{S}_i) = \sup_{f \in \mathcal{F}} \left(\sum_{i=1}^m f^2(\mathcal{S}_i) \right) - f^2(\mathcal{S}_j) .$$

Clearly we have that $0 \leq Z(\mathcal{S})$ and $\forall j \in 1 \dots m : Z(\mathcal{S}) - 1 \leq Z_j(\mathcal{S}) \leq Z(\mathcal{S})$, essentially because $\mathcal{F}^2(\mathcal{D}) \subseteq [0, 1]$. We now show

that $Z(\cdot)$ obeys $(1, 0)$ -self-boundedness.

$$\begin{aligned} \sum_{j=1}^m Z(\mathcal{S}) - Z_j(\mathcal{S}) &= mZ(\mathcal{S}) - \sum_{j=1}^m \sup_{f \in \mathcal{F}} \left(\left(\sum_{i=1}^m f^2(\mathcal{S}_i) \right) - f^2(\mathcal{S}_j) \right) \\ &\leq mZ(\mathcal{S}) - \sup_{f \in \mathcal{F}} \sum_{j=1}^m \left(\sum_{i=1}^m f^2(\mathcal{S}_i) \right) - f^2(\mathcal{S}_j) \\ &= mZ(\mathcal{S}) - \sup_{f \in \mathcal{F}} (m-1) \sum_{i=1}^m f^2(\mathcal{S}_i) \\ &= mZ(\mathcal{S}) - (m-1)Z(\mathcal{S}) = Z(\mathcal{S}) \end{aligned}$$

We may now conclude that $Z(\cdot)$ is a self-bounding function. It thus follows that $\beta = \frac{1}{m}Z(\mathcal{S})$ is self-bounding with scale $\frac{1}{m}$. The remainder of the result then follows via Thm. A.2. \square

THEOREM 3.1. *Let $\eta \in (0, 1)$, and define the quantities*

$$\begin{aligned} \gamma &\doteq \beta + \frac{2 \ln \frac{5}{\eta}}{3m} + \sqrt{\left(\frac{\ln \frac{5}{\eta}}{\sqrt{3m}} \right)^2 + \frac{2\beta \ln \frac{5}{\eta}}{m}}, \\ \rho &\doteq \hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda) + \frac{2 \ln \frac{5}{\eta}}{3km} + \sqrt{\frac{4\beta \ln \frac{5}{\eta}}{km}}, \\ r &\doteq \rho + \frac{\ln \frac{5}{\eta}}{3m} + \sqrt{\left(\frac{\ln \frac{5}{\eta}}{2\sqrt{3m}} \right)^2 + \frac{\rho \ln \frac{5}{\eta}}{m}}, \\ \varepsilon &\doteq 2r + \frac{\ln \frac{5}{\eta}}{3m} + \sqrt{\frac{2(\gamma + 4r) \ln \frac{5}{\eta}}{m}}, \end{aligned} \quad (7)$$

Then, with prob. at least $1 - \eta$ over the choice of \mathcal{S} and Λ , it holds

$$SD(\mathcal{F}, \mathcal{S}) \leq \varepsilon.$$

PROOF. We show this result as a union bound over *exponential tail bounds* (hence the $\ln \frac{5}{\eta}$ in all bounds) on five events, each dependent on the draw of Λ and \mathcal{S} . We now enumerate these bounds, and note that where each bound depends on an unknown quantity, we substitute an upper-bound based on an empirical quantity.

The first is on γ , which upper-bounds the wimpy-variance via Thm. A.3. The next is on ρ , which holds via the Monte-Carlo estimation theorem [24, Thm. 5]. The third is on r , and holds via the self-bounding function inequality (Thm. A.2), using the well-known scale $1/2m$ -self-bounding property of Rademacher averages [15, Ex. 3.12]. Finally, the fourth and fifth tail bounds are on ε (upper and lower tails), both via Bousquet's inequality [16], which bounds the SD as stated.

These five bounds then all hold simultaneously with probability $\geq 1 - 5 \cdot \frac{\eta}{5} = 1 - \eta$ via the *union bound*, yielding the desideratum. \square

We now prove Thm. 4.2, recalled below for convenience.

THEOREM 4.2. *Let $d \geq \text{vd}(G)$ or $d = 3$ if G is undirected and there is at most one SP between any pair of vertices in G . Let*

$$m^*(\bar{\epsilon}, \eta) \doteq \frac{1}{\bar{\epsilon}^2} \left(4Cd + 4\sqrt{\frac{Cd \ln \frac{2}{\eta}}{2}} + \frac{\ln \frac{2}{\eta}}{2} \right), \quad (10)$$

where C is a universal constant [33, 41, 60]. When run with inputs $m \geq m^*(\bar{\epsilon}, \eta)$, and $\delta = \eta$, BAVARIAN could return $(\tilde{B}, \bar{\epsilon})$ rather than using the computed ϵ , and Thm. 4.1 would still hold.

The Rademacher average of a family \mathcal{F} of functions w.r.t. the probability distribution π over the domain \mathcal{D} of the functions in \mathcal{F} is the expectation, w.r.t. both the sample \mathcal{S} (which is drawn according to π^m) and the set of m k -dimensional Rademacher vectors Λ of the k -MCERA:

$$R_m(\mathcal{F}, \pi) \doteq \mathbb{E}_{\mathcal{S}, \Lambda} \left[\hat{R}_m^k(\mathcal{F}, \mathcal{S}, \Lambda) \right].$$

The following result, a pinnacle of statistical learning theory, connects the Rademacher average to the supremum deviation

THEOREM A.4. *With probability at least $1 - \eta$ over the choice of \mathcal{S} , it holds*

$$SD(\mathcal{F}, \mathcal{S}) \leq 2R_m(\mathcal{F}, \pi) + \sqrt{\frac{\ln \frac{2}{\eta}}{2m}}. \quad (13)$$

In all cases of interest for us, \mathcal{F} and π are associated to a BC approximation method A (see Sect. 3.2), so we use the notation $R_m(A)$ to refer to $R_m(\mathcal{F}_A, \pi_A)$.

LEMMA A.5. *Let $d \geq \text{vd}(G)$ or $d = 3$ if G is undirected and there is at most one SP between any pair of vertices in G . For some universal constant C [33, 41, 60], it holds that*

$$R_m(rk) \leq \sqrt{\frac{Cd}{m}}.$$

PROOF. The thesis follows from the fact that d is an upper bound to the VC-dimension [63] of the BC estimation task on G , as shown by Riondato and Kornaropoulos [56, Thm. 2], and from a standard result from statistical learning theory connecting the Rademacher average to the VC-dimension [33, 60]. \square

The proof of Thm. 4.2 for the rk estimator then follows from using the upper bound in Lemma A.5 as an upper bound to the Rademacher average in Thm. A.4, requiring the l.h.s. of (13) to be at most $\bar{\epsilon}$, and solving for m to obtain the expression for $m^*(\bar{\epsilon}, \eta)$.

We show next that Thm. 4.2 holds also for the estimators ab and bp. To this end, we need to introduce the concepts of *unique-SP-enforcers* and *unique-SP-enforcing families*.

Definition A.6. A (*unique-SP*)-enforcer t is a function from $V \times V$ to the set of all SPs in the graph G , i.e.,

$$t : V \times V \rightarrow \bigcup_{(u,v) \in V \times V} \Gamma_{uv},$$

mapping each pair $(u, v) \in V \times V$ to a SP $p \in \Gamma_{uv}$ if such a path exists (i.e., if $\Gamma_{uv} \neq \emptyset$), and to \emptyset otherwise.

There are

$$q \doteq \prod_{\substack{(u,v) \in V \times V \\ \text{s.t. } \sigma_{uv} \geq 1}} \sigma_{uv} \quad (14)$$

different enforcers. Let Q denote the set of all enforcers. Clearly $q = |Q|$.

Definition A.7. Given an enforcer $t \in Q$, the *(unique-SP)-enforcing family* \mathcal{T}_t associated to t is a family of functions from to $\{0, 1\}$, such that \mathcal{T} contains one function $f_{t,w}$ for each $w \in V$ defined as

$$f_{t,w}(u, v) \doteq \begin{cases} 1 & \text{if } w \in \text{int}(t(u, v)) \\ 0 & \text{otherwise} \end{cases} .$$

The following result holds for every enforcing family \mathcal{T}_t , $t \in Q$.

LEMMA A.8. For any $t \in Q$, the VC-dimension of \mathcal{T}_t is at most $\text{vd}(G)$, unless $q = 1$, in which case the VC-dimension of the only \mathcal{T}_t is exactly 3.

PROOF. The proof for the general case is exactly the same as the one for [56, Thm. 2], while the one for the case $q = 1$ is the one for [56, Lemma 2]. \square

A family \mathcal{T}_t cannot generally be seen as the family resulting from imposing a set of weights on the edges of G enforcing unique shortest paths between every pair of connected vertices [20]. The reason is that t may be such that there exists $(u, v) \in V \times V$ such that there are two vertices (w, z) on the enforced SP $t(u, v)$ for which $t(w, z)$ is not a subset of $t(u, v)$. In other words, t does not enforce “consistent” SPs, which instead must happen when imposing a set of weights as above. For this reason, the result from [56, Lemma 2] can only be used when there is a single enforcer in Q , as it relies on such consistency of the SPs.

The following result connects the set $\{\mathcal{T}_t, t \in Q\}$ of enforcing families to the family \mathcal{F}_{ab} used by the ab estimator (see (3)).

LEMMA A.9. For any $w \in V$, let $f_w \in \mathcal{F}_{ab}$ be the function associated to w in this family, as defined in (3). For any $(u, v) \in \mathcal{D}_{ab}$, it holds

$$f_w(u, v) = \frac{1}{q} \sum_{t \in Q} f_{t,w}(u, v) .$$

PROOF. The thesis is immediate when $\sigma_{uv}(w) = 0$, as $f_{t,w}(u, v) = 0$ for each $t \in Q$. Let now $\Gamma_{uv}(w) \subseteq \Gamma_{uv}$ be the set of all SPs from u to v that w is internal to. From the definition of enforcer, it holds that for each SP $p \in \Gamma_{uv}(w)$, there are exactly

$$z \doteq \prod_{\substack{(x,y) \in V \times V \setminus \{(u,v)\} \\ \text{s.t. } \sigma_{xy} \geq 1}} \sigma_{xy} \quad (15)$$

enforcers in Q mapping (u, v) to p . Thus there are exactly $\sigma_{uv}(w) \cdot z$ enforcers that map (u, v) to a SP to which w is internal. It holds $f_{t,w}(u, v) = 1$ iff t is one of these $\sigma_{uv}(w) \cdot z$ enforcers. It follows that

$$\sum_{t \in Q} f_{t,w}(u, v) = \sigma_{uv}(w) \cdot z .$$

It also holds, from (15) and (14), that

$$\frac{1}{\sigma_{uv}} = \frac{z}{q} .$$

The thesis follows from these two identities as

$$\frac{\sigma_{uv}(w)}{\sigma_{uv}} = \frac{\sigma_{uv}(w)z}{q} = \frac{1}{q} \sum_{t \in Q} f_{t,w}(u, v) .$$

\square

The following lemma corresponds to Lemma A.5 for the ab estimator. By combining it with Thm. A.4, we obtain Thm. 4.2 for ab.

LEMMA A.10. Let $d \geq \text{vd}(G)$ or $d = 3$ if G is undirected and there is at most one SP between any pair of vertices in G . For some universal constant C [33, 41, 60], it holds that

$$R_m(ab) \leq \sqrt{\frac{Cd}{m}} .$$

PROOF. Let $\mathcal{S} = \{(u_1, v_1), \dots, (u_m, v_m)\}$ be a sample from \mathcal{D}_{ab} and $\Lambda = \{\lambda_1, \dots, \lambda_k\}$ be a bag of m -dimensional Rademacher vectors. From Lemma A.9 it holds

$$\hat{R}_k^m(\mathcal{F}_{ab}, \mathcal{S}, \Lambda) = \frac{1}{mk} \sum_{j=1}^k \sup_{w \in V} \sum_{i=1}^m \left(\lambda_{j,i} \frac{1}{q} \sum_{t \in Q} f_{t,w}(u_i, v_i) \right) .$$

From here, thanks to the fact that the constant zero function belongs to \mathcal{F} , we can use the subadditivity of the supremum operator to get

$$\begin{aligned} \hat{R}_k^m(\mathcal{F}_{ab}, \mathcal{S}, \Lambda) &\leq \frac{1}{q} \sum_{t \in Q} \frac{1}{mk} \sum_{j=1}^k \sup_{w \in V} \sum_{i=1}^m (\lambda_{j,i} f_{t,w}(u_i, v_i)) \\ &= \frac{1}{q} \sum_{t \in Q} \hat{R}_k^m(\mathcal{T}_t, \mathcal{S}, \Lambda) . \end{aligned}$$

The above inequality is true for every choice of \mathcal{S} and Λ , so it must also hold that

$$R_m(ab) \leq \frac{1}{q} \sum_{t \in Q} R_m(\mathcal{T}_t, \pi_{ab}) .$$

Lemma A.8 allows us to say that each summand on the r.h.s. is at most $\sqrt{(Cd)/m}$, and the proof is complete. \square

We now complete the discussion of Thm. 4.2 by showing results equivalent to Lemmas A.9 and A.10 for the bp estimator.

LEMMA A.11. For any $w \in V$, let $f_w \in \mathcal{F}_{bp}$ be the function associated to w in this family, as defined in (4). For any $u \in \mathcal{D}_{bp}$, it holds

$$f_w(u) = \frac{1}{q} \frac{1}{|V|-1} \sum_{v \neq w} \sum_{t \in Q} f_{t,w}(u, v) .$$

PROOF. The proof follows from the definition of $f_w(u)$ and of $f_w(u, v)$ in (3) and from Lemma A.9 \square

LEMMA A.12. Let $d \geq \text{vd}(G)$ or $d = 3$ if G is undirected and there is at most one SP between any pair of vertices in G . For some universal constant C [33, 41, 60], it holds that

$$R_m(bp) \leq \sqrt{\frac{Cd}{m}} .$$

PROOF. The proof follows the same step as the one for Lemma A.10, but using Lemma A.11 in place of Lemma A.9. \square

We now show the proof for Thm. 4.3.

THEOREM 4.3. With probability at least $1 - \delta$ (over the runs of the algorithm), the pair (\tilde{B}, ε) returned by BAVARIAN-P is such that

$$\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \varepsilon \leq \bar{\varepsilon} .$$

PROOF. Consider a variant V of BAVARIAN-P that, rather than sampling progressively, receives as input, in addition to all the parameters listed in Alg. 2, an *ordered sequence*

$$\mathcal{Z} = \{(s_1, \lambda_1), \dots, (s_{m_{T-1}}, \lambda_{m_{T-1}})\},$$

where $s_j \in \mathcal{D}_A$ is sampled independently (from anything else) according to π_A , and λ_j is an independently-sampled k -dimensional vector of independent Rademacher variables. The variant V works essentially as Alg. 2 with the exception that, instead of sampling elements of π_A and sampling vectors of Rademacher variables, it just picks the “next” element from \mathcal{Z} as needed.

The reason for introducing V is to make it evident that it is possible to reason about properties of \mathcal{Z} *independently* from the workings of the algorithm, as \mathcal{Z} is fixed before the algorithm even starts. Next, we show that V returns a pair (\tilde{B}, ε) with the desired properties mentioned in the thesis, and then we argue how to extend these properties to Alg. 2.

Fix $i \in [0, T - 1]$, and consider the value ε_i computed by V using the sequence of samples $\mathcal{Z}_i = \{(s_1, \lambda_1), \dots, (s_{m_i}, \lambda_{m_i})\}$. An application of Thm. 3.1 allows us to say that, with probability at least $1 - \delta/T$ (over the choice of \mathcal{Z}_i), it holds $\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \varepsilon_i$.

We can then apply the union bound over the T iterations, and obtain that, with probability at least $1 - \delta$ over its runs (i.e., over the choice of the sequence \mathcal{Z}), the output (\tilde{B}, ε) of V satisfies

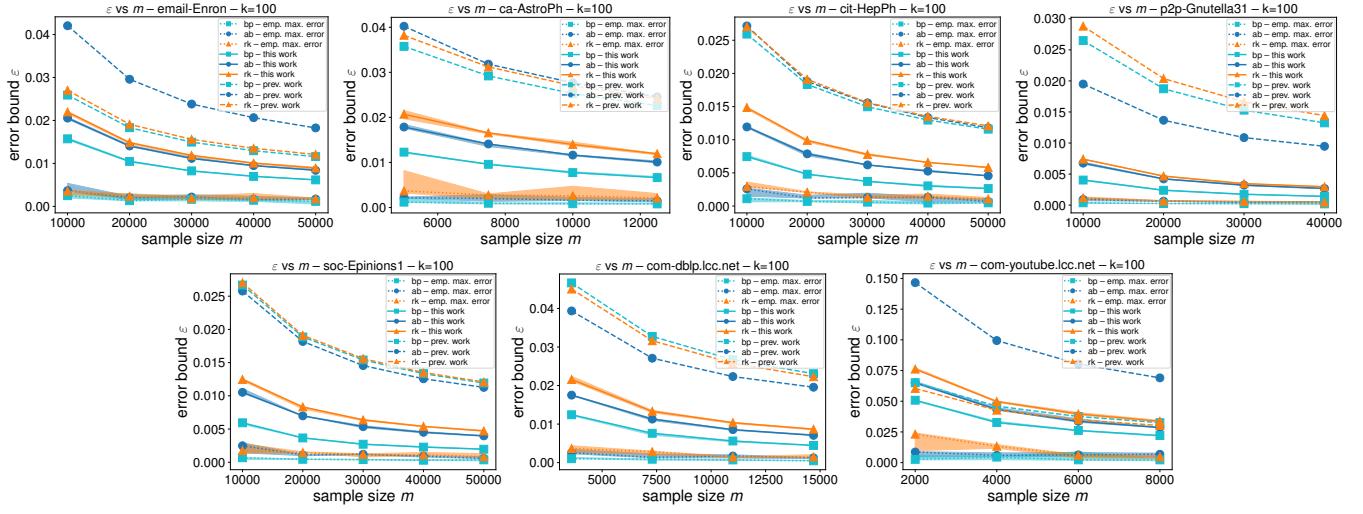
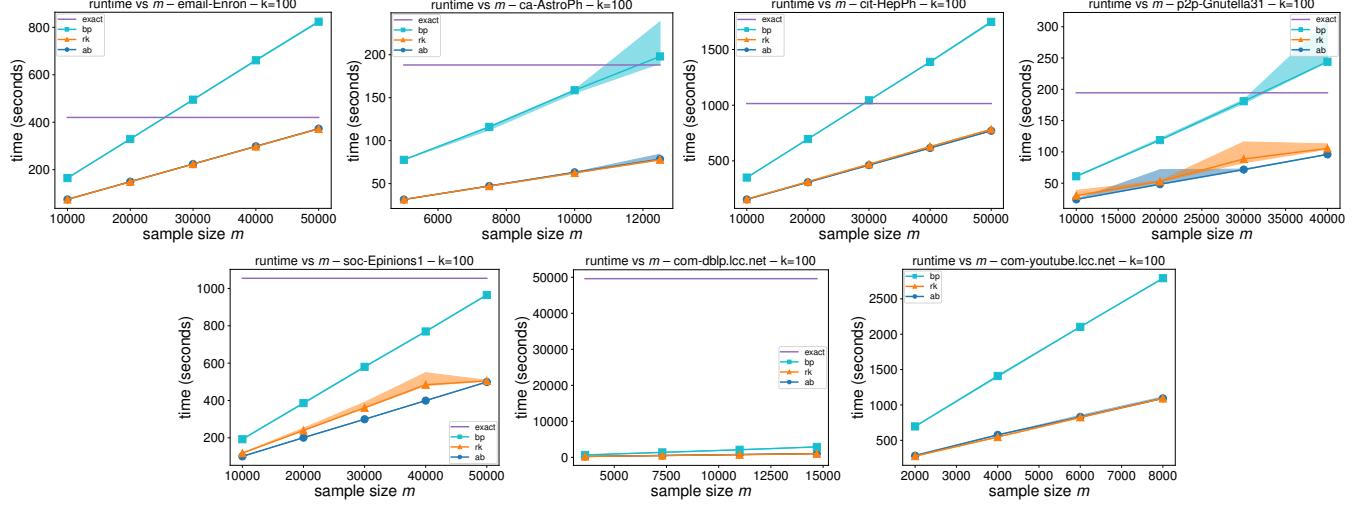
$$\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \varepsilon . \quad (16)$$

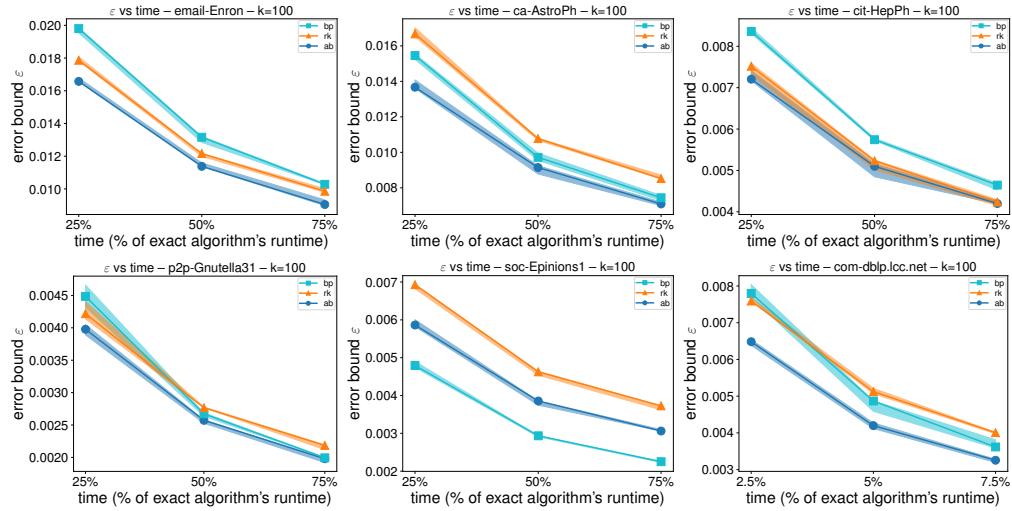
The last sample size m_{T-1} must satisfy the requirement from (11) and when using this sample size, V outputs $\varepsilon = \bar{\varepsilon}$. Combining this fact with (16), we get that the output (\tilde{B}, ε) of V satisfies

$$\max_{v \in V} |\tilde{b}_A(v) - b(v)| < \varepsilon \leq \bar{\varepsilon} .$$

We now argue that the properties enjoyed by V can be extended to BAVARIAN-P. Indeed to each execution of V with a set of parameters, we can associate an execution of BAVARIAN-P with the same set of parameters (minus \mathcal{Z}) such that the sequence of random bits used to generate \mathcal{Z} is given to BAVARIAN-P and used to draw samples from \mathcal{D}_A and the Rademacher vectors λ when needed. Clearly, the ordered sequence of samples and Rademacher vectors generated by BAVARIAN-P in this execution is a sub-sequence (potentially improper) of the sequence \mathcal{Z} given in input to V , and specifically exactly the sub-sequence of all and only the pairs that the execution of V “used” (the length of this sub-sequence is obviously m_i for some $i \in \{0, \dots, T - 1\}$). When using the sample size m_{T-1} , V returns $\bar{\varepsilon}$. The stopping condition on line 10 of Alg. 2 is guaranteed to be satisfied at some iteration, i.e., BAVARIAN-P will terminate. The pair (\tilde{B}, ε) returned by this execution of BAVARIAN-P is the same as the one returned by the associated execution of V , and this correspondence is true for all executions of the two algorithms. Hence, the same properties on the output of V also hold for BAVARIAN-P, and the proof is complete. \square

A.3 Additional Experimental Results

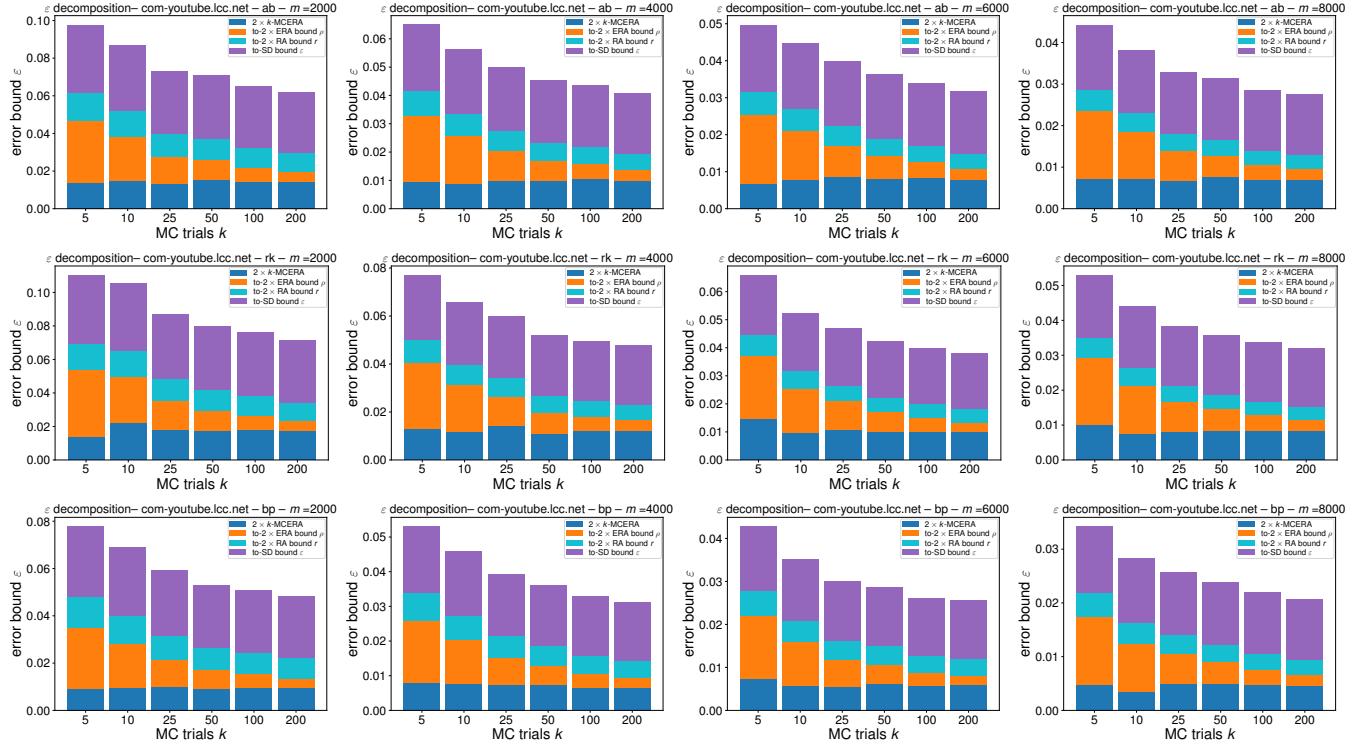
Figure 3: Error bound ε vs sample size m .Figure 4: Runtime vs sample size m .

Figure 5: Error bound ε vs time.







Figure 10: Error bound ε decomposition.