

EFFICIENT COMPUTATION OF THE WEIGHTED CLUSTERING COEFFICIENT

Silvio Lattanzi¹ and Stefano Leonardi²

¹Google Research, New York, New York, USA

²Università di Roma, “La Sapienza”, Roma, Italy

Abstract *The clustering coefficient of an unweighted network has been extensively used to quantify how tightly connected a neighbor is around a node, and it has been widely adopted for assessing the quality of nodes in a social network. The computation of the clustering coefficient is challenging because it requires counting the number of triangles in the graph. Several recent works have proposed efficient sampling, streaming, and MapReduce algorithms that make it possible to overcome this computational bottleneck.*

As a matter of fact, the intensity of the interaction between nodes, which is usually represented with weights on the edges of the graph, is also an important measure of the statistical cohesiveness of a network. Recently, various notions of weighted clustering coefficient have been proposed but all those techniques are hard to implement on large-scale graphs.

In this work we show how standard sampling techniques can be used to obtain efficient estimators for the most commonly used measures of weighted clustering coefficient. Furthermore, we propose a novel graph-theoretic notion of clustering coefficient in weighted networks.

1. INTRODUCTION

In recent years, we have observed a growing attention to the study of the structural properties of social networks [15, 17] as result of the fast increase of the amount of social network data available for research. A widely adopted measure of the graph structure of a social network is the clustering coefficient [33]. The local clustering coefficient of a node is defined as the probability that any two neighbors of a node are themselves neighbors. The clustering coefficient of a graph is the average local clustering coefficient of the nodes of the graph.

The clustering coefficient is used to measure how tightly interconnected the community is around a node. The degree of closeness of any two neighbors of a node is also interpreted as an index of trust of the node itself. The local clustering coefficient of a node has been proved, for example, to be a relevant feature for detecting spam nodes on the Web [3] and high-quality users in social networks [3].

Computing the clustering coefficient of a network is a challenging computational task because it reduces to counting the number of triangles in a graph. This task can be naïvely executed in $O(n^3)$ time or it can be reduced to matrix multiplication. The problem of computing the local clustering coefficient for every node of the network is even more challenging. Several recent works have proposed a variety of efficient methods for fast

Address correspondence to Stefano Leonardi, Università di Roma, “La Sapienza”, Via Ariosto 25, Roma 00185, Italy. Email: leonardi@diag.uniroma1.it

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/ujnm.

computation of clustering coefficient in large-scale networks based on random sampling [11], streaming algorithms [7, 13], and MapReduce parallel computation [29].

However, most of the studies on the structural properties of social networks have focused on unweighted networks. In practice, many real-world networks exhibit varying degrees of intensity and heterogeneity in the connections that are usually modeled with positive real weights on edges. Weights on edges are used, for instance, to measure the number of messages exchanged between friends or the number of links between hosts. Because the statistical level of cohesiveness in a network should, in principle, depend also on the weight of the edges, some recent interesting works have started to investigate weighted networks [21]. Several new notions of weighted clustering coefficient have also been introduced ([2, 23], among others), but, unfortunately, no efficient method for estimating the weighted clustering coefficient has been presented thus far.

Computing the exact values of the weighted clustering coefficient is at least as difficult as for the unweighted clustering coefficient. Sampling is the key for an efficient and accurate approximation [7, 11]. In the unweighted case, the key to the design of an unbiased estimator is the ability to draw uniformly at random a neighbor pair of a node and reporting 1 if and only if the neighbor pair is connected. The problem of drawing a neighbor pair can be efficiently solved in linear time if the two neighbors can be determined independently. The sampling of the two edges of a neighbor pair cannot be independent if the contribution to the clustering coefficient depends on the weights of the edges [2, 23], thus leading to a superlinear sampling complexity. Nevertheless, in this paper we show that, for several measures of weighted clustering coefficient, it is possible to obtain an efficient linear time estimator.

Our contributions. We summarize, following, the main contributions of our work:

1. We show how to obtain efficient estimators for several standard definitions of weighted clustering coefficient.
2. Our sampling algorithm are easily parallelizable too. We also develop a scalable MapReduce implementation of our estimators. Our implementation uses two rounds of MapReduce, it sends a number of messages across machines limited by the number of nodes times the number of samples required. The load for each machine is limited by the number of samples used by the algorithm times the maximum degree of a node in a graph.
3. We introduce a novel notion of *weighted clustering coefficient*. We base our proposal on the observation that edges with large weights are more likely to play a role in the social network. Our model defines a family of unweighted random graphs with edges existing with different probabilities. The probability of an edge depends on its weight. The larger the weight, the higher the probability. Each graph of the family of random graphs is an unweighted graph. The local weighted clustering coefficient of a node is defined as the expected local clustering coefficient in the family of random graphs. Our definition naturally extends to the weighted clustering coefficient of the entire graph.¹
4. We also design a polynomial time algorithm to compute the value of the weighted clustering coefficient and a sampling technique to estimate it efficiently. The computation of the weighted clustering coefficient in our model does not require the generation of

¹We note that our definition of weighted random graph is different from the definition of [10] and it is more in line with the standard definition used in data mining and biology [12].

all graphs of the family. That would be computationally prohibitive. We show that a dynamic program is able to compute the exact local weighted clustering coefficient in polynomial time.² The computational complexity of this exact computation is still prohibitive in practice. We are, however, able to design an efficient estimator also for this new definition of clustering coefficient.

5. We perform experiments that show interesting properties of the weighted clustering coefficient.

1.1. Related Works

A survey of several approaches to clustering coefficient in weighted networks can be found in [26]. In [24], the definition of clustering coefficient is based on the average weight on the edges of a triangle. In [2], the definition of the local clustering coefficient of a node depends only on the weights of the two edges incident to the node but not on the weight of the third edge of the triangle. In [20], it is adopted as the standard unweighted definition with the exception that triangles are weighted by the edge that closes the triangle. In [23], the weight is considered only in the numerator of the definition of clustering coefficient, whereas the denominator is the one of the unweighted case. In [36], the weight of a triangle is obtained by multiplying the weights of the edges. Other proposals that are substantially different from our approach can also be found in [14, 37]. The study of the clustering coefficient in several classes of random unweighted graphs can be found in [4].

The problem of estimating the clustering coefficient is closely related to the problem of counting the number of triangles in a graph. This is computationally expensive even on graphs of moderate size, because of the time complexity needed to enumerate all the length-two paths of the graph. Several works proposed efficient heuristics [16, 28] with computational results reported for graphs of large size. More recently, there are algorithms that have been designed under the MapReduce [9] programming model. Using a MapReduce infrastructure, [29] proposed algorithms for computing the exact number of triangles and the clustering coefficient of graphs. Randomized algorithms for counting triangles were also implemented under the MapReduce paradigm [25]. Finally, to estimate the total number of triangles in a graph, it is possible to use also matrix sketches [19], unfortunately, it is not clear how to extend this approach to the local clustering coefficient. A related measure is also the transitivity coefficient of a graph [22]. Techniques adopted for estimating the clustering coefficient usually extend to the transitivity coefficient.

A natural approach for problems in massive networks is to provide approximate solutions based on the application of data stream and random sampling algorithms. These algorithms usually provide an $(1 \pm \epsilon)$ approximation of the number of triangles with probability $1 - \delta$. The number of samples and amount of memory needed depends on the quality of the approximation. Data stream algorithms for estimating the number of triangles of a graph have been considered in [13, 32]. Semistreaming algorithms have been proposed in [3]. A sampling-based algorithm for estimating the clustering coefficient of a graph is given in [27].

²The problem of computing a core decomposition in an uncertain graph with different probabilities on edges has been considered in [5]. The authors show how to compute the expected degree of an uncertain graph in polynomial time. This method cannot, however, be applied to speed up the computation of the the novel notion of weighted clustering coefficient that we propose.

2. PRELIMINARIES

Let $G = (V, E)$ be an undirected graph with $n = |V|$ and $m = |E|$ edges. For every vertex $v \in V$, let $\mathcal{N}(v, G)$ denote its neighborhood, i.e., $\mathcal{N}(v, G) = \{u \in V : \exists(u, v) \in E\}$. The clustering coefficient $C_v(G)$ of a vertex $v \in V$ is defined as the probability that a random pair of its neighbors is connected by an edge, i.e., $C_v(G) := \frac{|\{(u, w) \in E : u, w \in \mathcal{N}(v, G)\}|}{\binom{|\mathcal{N}(v, G)|}{2}}$. In the case of $|\mathcal{N}(v, G)| < 2$, we define $C_v(G) := 0$. The clustering coefficient $C(G)$ of G is the average clustering coefficient of its vertices, i.e., $C(G) = \frac{1}{n} \cdot \sum_{v \in V} C_v(G)$.

Let us denote by $W(v, G) = \{(u, w) : u, w \in \mathcal{N}(v, G)\}$ the set of *wedges* of vertex v in graph G , i.e., the set of distinct paths of length two centered at v .

We denote by $w : E \rightarrow \mathbb{R}^+$ the positive weight on the edges of the graph. Let $W = \max_{e \in E} w(e)$ be the maximum weight of an edge. We normalize the edge weights in a way that their range varies in $[0, 1]$. Denote by $p : E \rightarrow [0, 1]$ the normalized weights. We denote with $\mathbf{1}_C$, the indicator variable for the event C . In the experimental section, we will use the following classic normalization: $p(e) = \frac{1}{1 + \log W/w(e)}$.

Finally, we say that we have an (ϵ, δ) estimator for a measure M , if we can estimate M within an ϵ multiplicative factor with probability at least $1 - \delta$.

2.1. Generalizations of Clustering Coefficient in Weighted Networks

In this article, we consider three generalizations of the clustering coefficient in weighted networks. In particular, we focus our attention to two definitions [2, 23] that well represent two general approaches to the problem: in one case the weights of the edges are added, in the other case they are multiplied. We additionally introduce a novel definition that is particularly relevant when the weights on the edges can be interpreted as probabilities.³

Onnela et al. The first definition of clustering coefficient that we consider was introduced by Onnela et al. [23]:

$$WC_v^{\text{Onnela}} = \frac{\sum_{(u, w) \in W(v, G)} \hat{w}(e(v, u)) \hat{w}(e(v, w)) \hat{w}(e(u, w))}{|\mathcal{N}(v, G)| (|\mathcal{N}(v, G)| - 1)},$$

where, with $w(e(v, u))$, we indicate the weight of the edge $e(v, u)$ and $\hat{w}(e(\cdot, \cdot)) = \frac{w(e(\cdot, \cdot))}{W}$.

Barrat et al. The second definition of clustering coefficient that we consider was introduced by Barrat et al. [2]:

$$WC_v^{\text{Barrat}} = \frac{\sum_{(u, w) \in W(v, G)} (w(e(v, u)) + w(e(v, w))) \mathbf{1}_{e(u, w)}}{(|\mathcal{N}(v, G)| - 1) (\sum_{v \in e} w(e))},$$

where $\mathbf{1}_{e(u, w)}$ is equal to 1 if the edge (u, w) exists and 0 otherwise.

Weighted clustering coefficient for probabilistic networks. The last measure that we analyze is novel. The basic idea is that the normalized weights can be interpreted as probabilities of existence of the edges in the graph. More formally, define the class of random graph $\mathcal{G}_{n, p}$ with edge e appearing independently with probability $p(e)$. Each graph $G' = (V, E') \in \mathcal{G}_{n, p}$ is an edge subset E' of E . The probability of G' is $p(G') = \prod_{e \in E'} p(e) \prod_{e \notin E'} (1 - p(e))$.

³This setting is particularly relevant when graphs are generated by using inference models [12].

The weighted clustering coefficient WC_v of a vertex $v \in V$ is defined as the expected clustering coefficient over the class of graphs $\mathcal{G}_{n,p}$: $WC_v^{\text{random}} = E_{G' \in \mathcal{G}_{n,p}} C_v(G')$.

3. COMPUTING THE WEIGHTED CLUSTERING COEFFICIENT IN PROBABILISTIC NETWORKS

In this section we give a polynomial algorithm to compute the new definition of weighted clustering coefficient efficiently. Note that at first sight our problem seems computationally very challenging because there are exponentially many possible realizations of the neighborhood of each node.⁴

Our first algorithmic contribution is to show that the problem is in P, we give an algorithm with complexity $O(|\mathcal{N}(v, G)|^4)$. Our algorithm is based on a dynamic program that computes incrementally the contribution of each neighbor pair to the clustering coefficient of each node.

Unfortunately, our exact algorithm is too slow to run on real networks where the maximum degree is typically very large (in the order of millions for Twitter or Google+); fortunately, in the next section, we show that the new measure has an efficient (ϵ, δ) estimator.

Recall that the unweighted clustering coefficient of a node v is defined as the probability that a randomly selected pair of its neighbors is connected by an edge; based on this, we can give an alternative definition of weighted clustering coefficient for probabilistic networks. Let $\chi(u, w)$ be a random variable that has value 1 if the randomly selected pair is (u, w) and 0 otherwise. We have $C_v(G) := \sum_{u,w \in \mathcal{N}(v,G) \wedge (u,w) \in E} Pr(\chi(u, w) = 1)$, where each pair is counted only once. In the following, we shorten $\mathcal{N}(v, G')$ to $\mathcal{N}'(v)$. Using this definition, we can rewrite the weighted clustering coefficient for v as $WC_v^{\text{random}} = E_{G' \in \mathcal{G}_{n,p}} [\sum_{u,w \in \mathcal{N}'(v) \wedge (u,w) \in E'} Pr(\chi(u, w) = 1 | G')]$.

Now, by defining $\xi(u, w)$ as a random variable that has value 1 if and only if $u, w \in \mathcal{N}'(v) \wedge (u, w) \in E'$, and by denoting with $\mathbf{1}_{\xi(u,w)}$ its indicator function, we have

$$\begin{aligned} WC_v^{\text{random}} &= E_{G' \in \mathcal{G}_{n,p}} \left[\sum_{u,w \in \mathcal{N}'(v) \wedge (u,w) \in E'} Pr(\chi(u, w) = 1 | G') \right] \\ &= E_{G' \in \mathcal{G}_{n,p}} \left[\sum_{u,w \in \mathcal{N}'(v)} \left(\mathbf{1}_{\xi(u,w)} Pr(\chi(u, w) = 1 | G') \right) \right] \\ &= \sum_{u,w \in \mathcal{N}'(v)} E_{G' \in \mathcal{G}_{n,p}} \left[\mathbf{1}_{\xi(u,w)} Pr(\chi(u, w) = 1 | G') \right] \\ &= \sum_{u,w \in \mathcal{N}'(v)} \left(Pr(\xi(u, w) = 1) * E_{G' \in \mathcal{G}_{n,p}} \left[\mathbf{1}_{\xi(u,w)} Pr(\chi(u, w) = 1 | G') \mid \xi(u, w) = 1 \right] \right) \\ &= \sum_{u,w \in \mathcal{N}'(v)} \left(Pr(u, w \in \mathcal{N}'(v) \wedge (u, w) \in E') * Pr(\chi(u, w) = 1 \mid \xi(u, w) = 1) \right). \end{aligned}$$

⁴Note that enumerating all the triangles in the graph would not work in this setting because of the dependency induced by the number of wedges in the realization of the random graph.

Now, the first term of the sum can be easily computed because $Pr(u, w \in \mathcal{N}'(v) \wedge (u, w) \in E') = p(e_{u,v})p(e_{w,v})p(e_{w,u})$. The second term is still problematic. In fact, $Pr(\chi(u, w) = 1 | \xi(u, w) = 1)$ depends on all the possible instantiations of G' , and so it potentially involves the computation of exponentially many terms.

In the following, we show how to compute it efficiently using dynamic programming.⁵ Note that $Pr(\chi(u, w) = 1 | \xi(u, w) = 1) = Pr(\chi(u, w) = 1 | u, w \in \mathcal{N}'(v))$ because the existence of the edge (u, w) does not change the probability of selecting u and w as random neighbors of v . And $Pr(\chi(u, w) = 1 | u, w \in \mathcal{N}'(v))$ is the probability that a pair u, w of neighbors of v are selected conditioned on the fact that $u, w \in \mathcal{N}'(v)$.

To compute this probability, we use the equivalence between the following two processes. The first selects two elements uniformly at random without replacement from a set S , and the second computes a random permutation of the elements in the set S and then returns the first two elements of the permutation.

Using this equivalence, we can rephrase the probability $Pr(\chi(u, w) = 1 | u, w \in \mathcal{N}'(v))$ as the probability that in a random permutation of the nodes in $\mathcal{N}(v)$, u , and w are the two nodes with the smallest positions in $\mathcal{N}'(v)$. Note that, for this to happen, either u and w are the first two nodes in the permutation of the nodes in $\mathcal{N}(v)$, or none of the nodes that are in positions smaller than u and w appear in $\mathcal{N}'(v)$.

Now, leveraging on this fact, we give a quadratic dynamic program to compute $Pr(\chi(u, w) = 1 | \xi(u, w) = 1)$. Consider an arbitrary order to the nodes in $\mathcal{N}(v) \setminus \{u, w\}$, $z_1, z_2, \dots, z_{|\mathcal{N}(v)|-2}$. In our algorithm, we implicitly construct all the possible permutations incrementally, and at the same time we estimate the probability that u, w are selected in each permutation. More specifically, initially we analyze the permutations containing only the elements $\{u, w\}$, then those containing the elements $\{u, w, z_1\}$, then those containing the elements $\{u, w, z_1, z_2\}$, and so on so until we get the probability for each permutation containing all the elements in $\mathcal{N}(v)$.

The key ingredient of our algorithm is the following observation. Once we have computed the probability for all the permutations containing the nodes $\{u, w, z_1, z_2, \dots, z_{i-1}\}$, to extend our computation to the permutations containing also the node z_i , we have to consider only two scenarios: in the first z_i appears after u, w in the permutation; in this case, the probability that u and w are the nodes in $\mathcal{N}'(v)$ with the two smallest positions is the same. In the second, z_i appears before either of u or of w , conditioned on this event, the probability that u and w are the nodes in $\mathcal{N}'(v)$ with the two smallest positions decreases by a multiplicative factor $1 - p(e_{v,z_i})$.

We are now ready to state our dynamic program more formally. Let M be a square matrix of dimension $|\mathcal{N}(v)| - 1$ such that $M_{i,j}$, for $j \leq i$, contains the probability that in a random permutation of nodes $\{u, w, z_1, z_2, \dots, z_i\}$ u and w are preceded by exactly j elements in the permutation but they are in the first and second positions when we consider the ordering induced only to nodes in $\mathcal{N}'(v)$. Note that $M_{0,0}$ is equal to 1, because in this case, we consider permutations containing only $\{u, w\}$. Similarly, we can compute $M_{1,0}$ and $M_{1,1}$. In particular, for $M_{1,0}$ we require that z_1 is in a position after u and w , so we have $M_{1,0} = \frac{1}{3}M_{0,0}$. Instead, $M_{1,1} = \frac{2}{3}(1 - p(e_{v,z_1}))M_{0,0}$. More generally, we have that for

⁵Unfortunately, to the best of our knowledge, there is no analytic technique to estimate this quantity correctly or with a close approximation without using a dynamic programming.

$j \leq i$,

$$M_{i,j} = \begin{cases} \frac{i-1}{i+1} M_{i-1,0} & \text{if } j=0 \\ \frac{i-j-1}{i+1} M_{i-1,j} + \frac{j+1}{i+1} \bar{p}(e_{v,z_i}) M_{i-1,j-1} & \text{if } j < i \\ & \text{and } j > 0 \\ \frac{i}{i+1} \bar{p}(e_{v,z_i}) M_{i-1,j-1} & \text{if } j = i \end{cases},$$

where $\bar{p}(\ast) = 1 - p(\ast)$.

Once we have computed the matrix M , we can compute $Pr(\chi(u, w) = 1 | u, w \in \mathcal{N}'(v))$, in fact, we have that $Pr(\chi(u, w) = 1 | u, w \in \mathcal{N}'(v)) = \sum_{i=0}^{|\mathcal{N}'(v)|-2} M_{|\mathcal{N}'(v)|-2,i}$. So we have

$$WC_v^{\text{random}} = \sum_{u,w \in \mathcal{N}(v)} \left(\frac{1}{2} p(e_{u,v}) p(e_{w,v}) p(e_{w,u}) \ast \left(\sum_{i=0}^{|\mathcal{N}'(v)|-2} M_{|\mathcal{N}'(v)|-2,i} \right) \right).$$

We summarize our algorithm to compute WC_v^{random} here:

Algorithm (exact WC_v^{random})

Input: The weighted subgraph induced by $v \cup \mathcal{N}(v)$.

Output: WC_v^{random} .

$WC_v^{\text{random}} = 0$.

for all $u, w \in \mathcal{N}(v)$ **do**

Compute the matrix M for u, w

Using M , compute the probability p that (u, v, w) is a triangle and is selected

$WC_v^{\text{random}} += p$

Output WC_v^{random} .

Note that the above algorithm has complexity $O(|\mathcal{N}(v)|^4)$, so it is too slow to run on large networks. For this reason, we study in the next section efficient estimators for the weighted clustering coefficient.

4. EFFICIENT ESTIMATORS FOR THE WEIGHTED CLUSTERING COEFFICIENT

We propose efficient (ϵ, δ) estimators for the various definitions of weighted clustering coefficient. Our estimators that use basic concentration theory are similar to that presented in [6, 31]. They are the first linear estimators for the weighted clustering coefficient, to the best of our knowledge.

Onnela et al. Recall the definition of Onnela et al. [23] given in Section 2.

In this definition, the weighted clustering coefficient is equal to the total normalized weight of the triangles containing v , averaged by the number of wedges centered on v . Thus,

if we sample the wedges uniformly at random, using the Hoeffding bound and the fact that the normalized weights⁶ are in $[0, 1]$, we get an efficient (ϵ, δ) estimator for WC_v^{Onnela} .

More formally, let X_1, \dots, X_s be identical random variables that, with probability $\frac{1}{|\mathcal{N}(v, G)|(|\mathcal{N}(v, G)|-1)}$, have value $\hat{w}(e(v, u))\hat{w}(e(v, w))\hat{w}(e(u, w))$ for every wedge $\langle u, w \rangle$. Then, $E[\sum_{i=1}^s X_i] = sWC_v^{\text{Onnela}}$. Furthermore, by the Hoeffding bound, we have that: $P(|X - E[\sum_{i=1}^s X_i]| \leq \epsilon E[\sum_{i=1}^s X_i]) \leq e^{-\frac{\epsilon^2 E[\sum_{i=1}^s X_i]}{3}} = e^{-\frac{\epsilon s WC_v^{\text{Onnela}}}{3}}$. So, if we want $\delta > e^{-\frac{\epsilon s WC_v^{\text{Onnela}}}{3}}$, it suffices to have $s \in O(\log \frac{1}{\delta} \cdot \frac{1}{\epsilon^2 \cdot WC_v^{\text{Onnela}}})$ samples.

Lemma 4.1. *There is a sampling-based algorithm that, with probability $1 - \delta$, returns a $(1 \pm \epsilon)$ -approximation of the local weighted clustering coefficient WC_v^{Onnela} of a vertex v of a weighted graph G . It needs $O(\log \frac{1}{\delta} \cdot \frac{1}{\epsilon^2 \cdot WC_v^{\text{Onnela}}})$ samples.*

In the following, we present the pseudocode to compute this estimator:

Algorithm (sampling WC_v^{Onnela})

Input: The weighted subgraph induced by $v \cup \mathcal{N}(v)$.
Output: Approximate WC_v^{Onnela} .
for all $i = 1$ **to** s **do**
 sample a random wedge $\langle u, w \rangle$ **uniformly from** $\mathcal{N}(v)$
 If $(u, w) \in E$ **then set** $X_i \leftarrow w(u, v)^{1/3}w(u, w)^{1/3}w(v, w)^{1/3}$
 else set $X_i \leftarrow 0$
Output $X := \frac{1}{s} \cdot \sum_{i=1}^s X_i$.

Furthermore, note that, for the sampler, we need to be able to sample only random wedges and this can be easily done in linear time.

Barrat et al. Recall the definition of Barrat et al. [2] given in Section 2. In this case, the weighted clustering coefficient is not an explicit average, so we cannot use the Hoeffding bound directly as before. Nevertheless, note that we can define WC_v^{Barrat} as the average value of the random variable X , where X has value $\mathbf{1}_{e(u, w)}$ with probability $\frac{w(e(v, u)) + w(e(v, w))}{(|\mathcal{N}(v, G)|-1)(\sum_{v \in e} w(e))}$ for all $\langle u, w \rangle \in W(v, G)$.

Using this alternative definition combined with the Chernoff bound, we get that by using k samples of the wedges weighted with the correct probability we can get good estimation of WC_v^{Barrat} .

More formally, let X_1, \dots, X_s identical random variable that with probability $\sum_{\langle u, w \rangle: \mathbf{1}_{e(u, w)}=1} \frac{(w(e(v, u)) + w(e(v, w)))}{(|\mathcal{N}(v, G)|-1)(\sum_{v \in e} w(e))}$ have value 1 or 0 otherwise. Then, $E[\sum_{i=1}^s X_i] = sWC_v^{\text{Barrat}}$. Furthermore, by Chernoff bound we have that

$$P\left(\left|X - E\left[\sum_{i=1}^s X_i\right]\right| \leq \epsilon E\left[\sum_{i=1}^s X_i\right]\right) \leq e^{-\frac{\epsilon^2 E[\sum_{i=1}^s X_i]}{3}} = e^{-\frac{\epsilon s WC_v^{\text{Barrat}}}{3}}.$$

So, if we want $\delta > e^{-\frac{\epsilon s WC_v^{\text{Barrat}}}{3}}$, it suffices to have $s \in O(\log \frac{1}{\delta} \cdot \frac{1}{\epsilon^2 \cdot WC_v^{\text{Barrat}}})$ samples.

⁶Note that, for this to work, it is fundamental that the weight on the edges has been normalized, and so are in $[0, 1]$.

Lemma 4.2. *There is a sampling-based algorithm, which, with probability $1 - \delta$, returns a $(1 \pm \epsilon)$ -approximation of the local weighted clustering coefficient WC_v^{Barrat} of a vertex v of a weighted graph G . It needs $\mathcal{O}(\log \frac{1}{\delta} \cdot \frac{1}{\epsilon^2 \cdot WC_v^{\text{Barrat}}})$ samples.*

At first sight, it might look as though we need quadratic time to sample a wedge with the correct probability, but also, in this case, it is possible to get a sample in linear time. In fact, to sample an edge with the correct probability, it is enough to sample the first edge e_1 with probability $p_1(e_1) = \frac{(|\mathcal{N}(v,G)|-2)w(e_1)+(\sum_{v \in e} w(e))}{2(|\mathcal{N}(v,G)|-1)(\sum_{v \in e} w(e))}$ and the second one e_2 with probability $p_2(e_2|e_1) = \frac{w(e_1)+w(e_2)}{(|\mathcal{N}(v,G)|-2)w(e_1)+(\sum_{v \in e} w(e))}$. The probability to sample pair e_1, e_2 is exactly $p_1(e_1)p_2(e_2|e_1) + p_1(e_2)p_2(e_1|e_2) = \frac{w(e_1)+w(e_2)}{(|\mathcal{N}(v,G)|-1)\sum_{v \in e} w(e)}$. However, it is also easy to verify that $\sum_e p_1(e) = 1$ and $\sum_{e \neq e_1} p_2(e|e_1) = 1$. We present the pseudocode for this estimator also:

Algorithm (sampling WC_v^{Barrat})

Input: The weighted subgraph induced by $v \cup \mathcal{N}(v)$.

Output: Approximate WC_v^{Barrat} .

for all $i = 1$ **to** s **do**

sample the first edge $e_1 = (u, v)$ **of the wedge with probability equal to** $p_1(e_1) = \frac{(|\mathcal{N}(v,G)|-2)w(e_1)+(\sum_{v \in e} w(e))}{2(|\mathcal{N}(v,G)|-1)(\sum_{v \in e} w(e))}$

sample the second edge $e_2 = (u, w)$ **of the wedge with probability equal to** $p_2(e_2|e_1) = \frac{w(e_1)+w(e_2)}{(|\mathcal{N}(v,G)|-2)w(e_1)+(\sum_{v \in e} w(e))}$

If $(u, w) \in E_i$ **then set** $X_i \leftarrow \frac{w(u,v)+w(v,w)}{(|\mathcal{N}(v,G)|-1)\sum_{v \in e} w(e)}$

else set $X_i \leftarrow 0$

Output $X := \frac{1}{s} \cdot \sum_{i=1}^s X_i$.

Weighted clustering coefficient for probabilistic networks. The algorithm is based on sampling a random wedge $\langle u, w \rangle \in W(v, G')$ from a random graph $G' \in \mathcal{G}_{n,p}$ and checking whether $(u, w) \in G'$. The core idea of our sampler is to generate for a node v , s neighbor realizations $\mathcal{N}(v)_1, \dots, \mathcal{N}(v)_s$ uniformly at random from $\mathcal{G}_{n,p}$. Then, for each realization, sample a random wedge $\langle u, w \rangle$ uniformly from $\mathcal{N}(v)_i$ and check if the wedge is part of a triangle in the realization. The estimation of the clustering coefficient is equal to the number of wedges that are part of a triangle divided by s .

Algorithm (sampling WC_v^{random})

Input: The weighted subgraph induced by $v \cup \mathcal{N}(v)$.

Output: Approximate WC_v^{random} .

Sample s neighbor realization $\mathcal{N}(v)_1, \dots, \mathcal{N}(v)_s$ uniformly at random from $\mathcal{G}_{n,p}$

for all $i = 1$ **to** s **do**

sample a random wedge $\langle u, w \rangle$ **uniformly from** $\mathcal{N}(v)_i$

If $(u, w) \in E_i$ **then set** $X_i \leftarrow 1$

else set $X_i \leftarrow 0$

Output $X := \frac{1}{s} \cdot \sum_{i=1}^s X_i$.

For the sake of completeness, we give a simple analysis of the algorithm, following. We first show that the expected value of X_i is exactly WC_v^{random} .

We have for each $i \in \{1, \dots, s\}$: $\mathbf{E}[X_i] = \mathbf{E}_{G' \in \mathcal{G}_{n,p}} \left[\frac{|\{(u,w) \in E': u,w \in \mathcal{N}(v,G')\}|}{\binom{\mathcal{N}(v,G')}{2}} \right] = \mathbf{E}_{G' \in \mathcal{G}_{n,p}} [C_v(G')] = WC_v^{\text{random}}$.

Then, we use the fact that, for $0-1$ random variables, we have $\mathbf{Var}[X_i] \leq \mathbf{E}[X_i^2] = \mathbf{E}[X_i] = WC_v^{\text{random}}$.

Now, we analyze the variance of X . Since the X_i are mutually independent, we get $\mathbf{Var}[X] = \mathbf{Var}\left[\frac{1}{s} \cdot \sum_{i=1}^s X_i\right] = \frac{1}{s^2} \cdot \sum_{i=1}^s \mathbf{Var}[X_i] \leq \frac{WC_v^{\text{random}}}{s}$. Finally, we can apply the Chebyshev inequality. This gives us $\Pr[|X - \mathbf{E}[X]| \geq \epsilon \cdot \mathbf{E}[X]] \leq \frac{\mathbf{Var}[X]}{(\epsilon \cdot \mathbf{E}[X])^2} \leq \frac{WC_v^{\text{random}}}{s \cdot \epsilon^2 \cdot (WC_v^{\text{random}})^2} = \frac{1}{s \cdot \epsilon^2 \cdot WC_v^{\text{random}}}$.

If $s \geq \frac{3}{\epsilon^2 \cdot WC_v^{\text{random}}}$, then with probability $\frac{2}{3}$, the algorithm sampling WC_v^{random} approximates the weighted clustering coefficient of vertex v in G within a relative error of $(1 \pm \epsilon)$. In order to amplify the probability of success, we run the algorithm $\Theta(\log \frac{1}{\delta})$ times and return the median of all results. This leads to the following corollary:

Lemma 4.3. *There is a sampling-based algorithm, which with probability $1 - \delta$, returns a $(1 \pm \epsilon)$ -approximation on the local weighted clustering coefficient WC_v^{random} of a vertex v of a weighted graph G . It needs $\mathcal{O}(\log \frac{1}{\delta} \cdot \frac{1}{\epsilon^2 \cdot WC_v^{\text{random}}})$ samples.*

5. PARALLEL IMPLEMENTATION

In this section, we first give a brief introduction to the MapReduce [9] framework and then we describe a highly optimized and scalable MapReduce implementation of our sampling algorithms. In particular, here we focus on parallelizing the algorithm to estimate WC_v^{random} ; the other sampling strategies can be parallelized in a very similar way. In the experimental section, we report experimental results showing that this implementation is extremely fast in practice.

The MapReduce framework is designed to simplify the implementation of parallel algorithms at very large scale. In the MapReduce framework, the data is processed in tuples composed of $\langle key, value \rangle$. The computation proceeds in rounds. In the Map phase, each machine receives all the values associated with a specific key k , then it executes some computation and output $\langle key, value \rangle$ tuples with potentially different key k' . A Shuffle phase aggregates all tuples with same key k' that are sent to the same physical machine. Finally, in the Reduce phase, each machine performs a computation that depends only on the tuples with the same key k' outputted from the Mapper, and output $\langle key, value \rangle$ tuples with key equal to the input k' .

To write a MapReduce program, it is typically important to design an algorithm that (i) minimizes the number of MapReduce rounds that are involved; (ii) minimizes the amount of communication between machines; and (iii) balances the working load across different machines. In the following, we show how these requirements are achieved in the implementation of our sampling algorithm in MapReduce.

We assume that the input graph is stored in $\langle key, value \rangle$ tuples, representing the adjacency list of each node. In the first Map phase, each machine reads the adjacency list of a node u . For sample $i = 1, \dots, s$, the machine constructs a realization of the neighborhood of a node u , $\mathcal{N}_i(u)$, according to $\mathcal{G}_{n,p}$ and samples a pair of random neighbors $(v_i, w_i) \in \mathcal{N}_i(u)$. Then it sends a message with key w_i and value $i, (u, v_i)$ to the machine

that controls node w_i .⁷ The informal meaning of these messages is that node u asks node w_i whether edge (w_i, v_i) exists in the i -th realization so that we can infer that triangle u, v_i, w_i exists in realization i . Finally, node u also sends its adjacency list to itself in order to be able to answer the requests of other nodes.

In the first Reduce phase, node u receives its own adjacency list and various requests $i, (w, v_i)$ to check the existence of edge (u, v_i) in realization $\mathcal{N}_i(u)$. If the test is positive, it writes a value $\langle u, w \rangle$ with its own key indicating that it is incident to a triangle with node w in one of the samples.

In the second Map phase, each node v reads the values written in the previous Reduce phase and, for each detected triangle $\langle v, u \rangle$, sends a message $\langle u, 1 \rangle$ to the other node u certifying the existence of the triangle. Finally, in the last Reduce step, each node receives the number of sampled triangles and simply computes its clustering coefficient by dividing it by the number of samples.

The implementation presented above uses two rounds of MapReduce, it sends a number of messages across machines upper bounded by the number of nodes times the number of samples required. The load for each machine is upper bounded by the number of samples used by the algorithm times the maximum degree of a node in a graph.

6. EXPERIMENTS

The main goal of this section is to show experimentally some properties of the weighted clustering coefficient and to show the speed-up obtained by our simple estimators.

We start by describing a classical application of the clustering coefficient and the dataset that we will use in our experiment.

Then, we analyze experimentally different techniques to map integer weights to weights between $[0, 1]$. In particular, we compare two mapping techniques and we analyze the trade-off between them. Note that this mapping is needed for two of our three clustering coefficient measures.⁸

We then compare the performance of the three weighted clustering coefficient measures with the classic unweighted clustering. Our findings are quite encouraging, in fact, we observe that the weighted clustering performance is always at least as good as the unweighted clustering coefficient and our new notion of weighted clustering coefficient outperforms the classic unweighted notion.

Finally, we analyze the scalability of our approach. In particular, we run our algorithm using a different number of machines on networks of increasing sizes. We observe that our algorithm is highly scalable and it can fully leverage on parallelization to improve its performances.

Dataset and experiment settings. The clustering coefficient is a fundamental topological property of networks and also one of the most used topological features in

⁷Note that a naïve implementation of the sampling procedure would have running time quadratic in the size of the adjacency list. Fortunately, this is not necessary; in fact, it is possible to select a random pair of neighbors in linear time. In particular, it is enough to assign to each neighbor a random number and then select the two neighbors with the smallest assigned values. In this way, each pair of nodes has the same probability of being selected, and so we can obtain a random sample.

⁸Experimentally, we observed that a nonlinear mapping performs better than the linear mapping proposed by Onnela et al.

Dataset	Nodes	Edges
Patents[18]	3,774,768	16,518,948
Pocec[30]	1,632,803	30,622,564
LiveJournal[1]	4,847,571	68,993,773
Orkut[34]	3,072,441	117,185,083
Friendster[34]	65,608,366	1,806,067,135

Table I Network statistics.

machine learning on graphs. Indeed, it has been used to detect spam on the Web [3] and malicious users in social networks [35].

For this reason, we study the effectiveness of weighted clustering coefficient by studying its power as a machine learning feature. In particular, we focus on the specific case in which we are interested in detecting spam on the Web. Toward this end, we use a publicly available dataset [8] composed of a collection of hosts manually labeled (spam/nonspam) by a group of volunteers and by the weighted host graph network. The graph is composed of 114,529 hosts in the .UK domain, and there are 5709 hosts marked as “nonspam” and 344 hosts marked as “spam.” Even if the web graph is directed in this section, we ignore the directionality of the edges for simplicity.⁹ Finally, we note that there are 2058 hosts marked as nonspam and 93 hosts marked as spam with clustering coefficient larger than 0 (for any weighted or unweighted definition of clustering coefficient).

In our experiments we are interested in analyzing only the correlation between various definitions of the clustering coefficient and the integrity of a host. To do it, for each definition, we first compute the corresponding score for each labeled node, then we rank all the labeled nodes with scores larger than 0 according to their scores and compute the precision of each position i of the ranking as the percentage of nonspam hosts before position i . This measure, even if simplistic, gives a good intuition of the correlation between the clustering coefficient and the goodness of a page.

Finally, to analyze the scalability of our algorithm, we consider five graphs from the Stanford Network Analysis Project (SNAP) repository: Friendster, Orkut, LiveJournal, Pocec, and Patents. Those graphs are unweighted, so we assign a random weight between $[0, 1]$ to every edge, independently. In Table I we report some basic statistics on the graphs.

6.1. Building the Probabilistic Graph

Two of the analyzed definitions of weighted clustering coefficient cannot be applied if the edges’ weights are not in $[0, 1]$. In this subsection we analyze different techniques to build a probabilistic graph from the input graph by mapping the weights to probability in $[0, 1]$. As a case of study, we analyze the effect of different mappings on the .UK domain graph.

Let us define e_W and e_w , respectively, the maximum and the minimum weight of an edge in the input graph. A first natural technique to construct our probabilistic graph is to use

⁹Note that all the discussed notion of the clustering coefficient can be extended to capture the directionality of the edges.

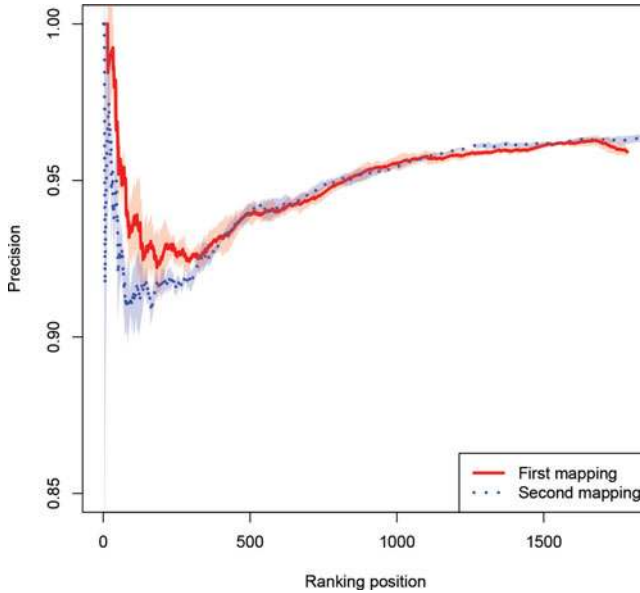


Figure 1 Comparison between the two mapping strategies M_1 and M_2 .

a linear mapping between $[e_w, e_W]$ and $[0, 1]$. This mapping, though, has a serious drawback for our probabilistic definition: in practice, the weights on the edges are distributed as a power law and so e_w/e_W is very small (for example, in our case of study, it is $1/2579857$). So, if we use this mapping, we would have very small weights on the edges of the graph and this would, in turn, imply an extremely small realization probability for every triangle in the graph.

To solve this issue in our experiment, we consider two nonlinear mapping functions M_1, M_2 . Both functions are mapping between $[e_w, e_W]$ and $[0, 1]$; more formally, we have that both $M_1, M_2 : [e_w, e_W] \rightarrow [0, 1]$. In particular, we define $M_1(w) = \frac{\log(w-e_w+1)}{\log(e_W-e_w+1)}$ and $M_2(w) = \frac{1}{1+\log(\frac{e_W-e_w+1}{w-e_w+1})}$. To compare those two mappings, we run our approximation algorithm for estimating WC_v^{random} for all the nodes in the graph, and we compare the precision of rankings that we obtain by using the two different rankings. In this experiment, to compute the weighted clustering coefficient, we execute 3200 samples per node, and to compute the average precision and the standard deviation we rerun the algorithm four times with different random seeds. In Figure 1 and in the rest of this section, we plot the average precisions using lines and the standard deviations using shadows around the lines.

From the experiments, it is possible to conclude that the two mapping strategies have similar performances, although M_1 seems to perform slightly better. For this reason, for the rest of this section we focus only on results obtained using the mapping M_1 .

6.2. Performances of the Sampling Algorithm

Now that we have defined our mapping strategy and built our probabilistic graph, we can focus on the performance of our sampling algorithm. Here, we first analyze the running time of the sampling algorithm presented in Section 4 when we vary the number of

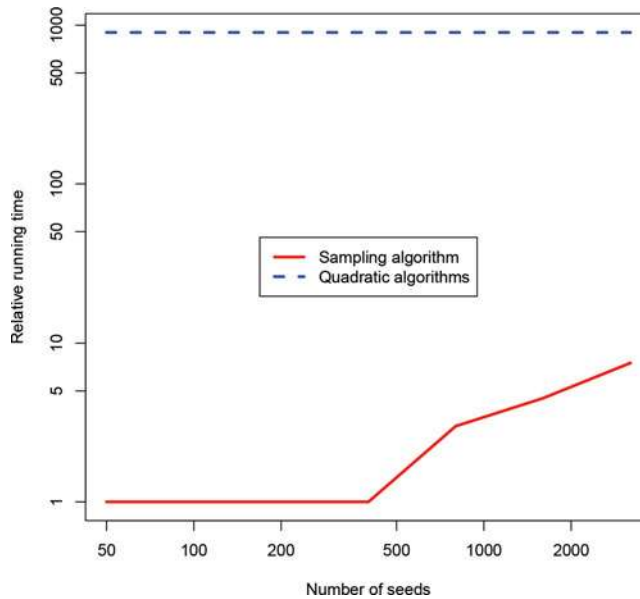


Figure 2 Running time vs. number of seeds and precision when we vary the number of samples between 50, 100, 200, and 3200.

samples used in the algorithm and we compare its running time with the running time of the algorithms that consider all the triangles to compute the unweighted clustering coefficient or the weighted clustering coefficient defined by Barrat et al. [2].

Then, we analyze how the precision of the ranking varies as a function of the number of samples performed by the algorithm. In Figure 2 we show the average running time of the sampling algorithm when we vary the number of samples relative to the running time of the optimal algorithm. It is interesting to note that the running time increases almost linearly with the number of seeds, showing that the algorithm efficiently uses all the parallelization offered by the MapReduce framework. Furthermore, it is quite interesting to note the huge difference in running time between the sampling algorithm and the quadratic algorithm that considers all the triangles. In fact, when we used 50, 100, 200, and 400 samples, the sampling algorithm was 900 times faster than the quadratic algorithm, and even when we used 2000 samples, the sampling algorithm is still 120 times faster!

Now, we turn our attention to the effects of varying the number of samples on the precision of the algorithm. In Figure 3 we show how the precision curve of the new notion of weighted clustering coefficient changes when we use 50, 100, 200, or 3200 samples (we notice a similar trend also with 400, 800, 1600 samples and for other clustering coefficient definitions, we do not show them in the figure, in order to maintain clarity). Also, in this case, we plot the average precisions with lines and the standard deviations with the shadows around the lines. From the plots, it seems that few samples are enough to obtain a good estimation of the weighted clustering coefficient.

There are several interesting observations to make here.

First, as predicted from Lemma 4.3, for all the measures the standard deviation decreases very quickly as the number of samples used by our algorithm increases.

Second, for WC_v^{random} the length of the ranking computed by our algorithm decreases when we use a smaller number of samples. This is probably due to the fact that several

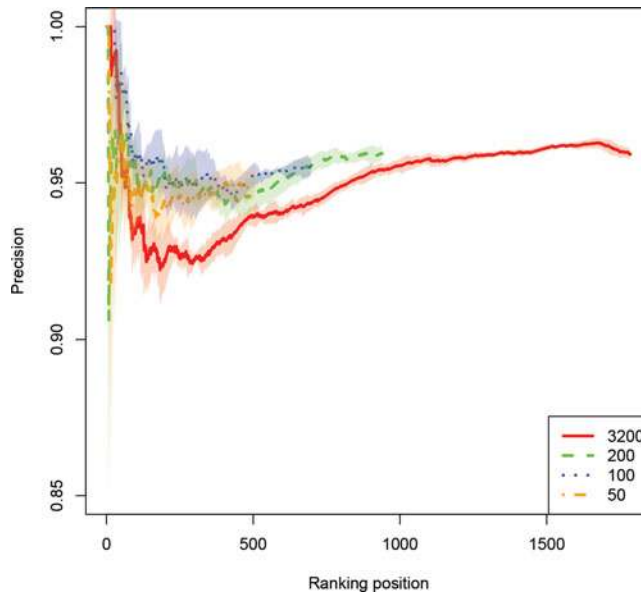


Figure 3 Precision vs. ranking position of the sample algorithm when we vary the number of samples among 50, 100, 200, and 3200.

nodes are incident to a small number of triangles and so by executing a small number of samples we do not discover them.

The third observation is probably the most striking: the precision of ranking provided by WC_v^{random} decreases when we use larger numbers of samples. We hypothesize that this phenomenon can be explained using the same explanation that we used for our second observation. In fact, also in this case, nodes that have small degrees are not likely to appear in the ranking when we consider few samples. But for nodes of small degrees, the clustering coefficient is probably not a meaningful indicator of their trustfulness. To verify this hypothesis, in the next subsection we analyze how the precision of the rankings changes when we consider only nodes with degrees above a specified threshold.

Motivated from the last observation, here we analyze the relationship between the degree of a node and the correlation between its clustering coefficient and its trustfulness for WC_v^{random} . To do this, we analyze the precision of the rankings of nodes when we are restricted only to nodes with weighted or unweighted degree above a specific threshold.

In Figure 4, we analyze the precision of the rankings computed by our sampling algorithm by using 3200 samples when we restrict computation to nodes with unweighted degree at least 0, 5, 10, and 20.

We observe a trend similar to that observed in Figure 3, suggesting that there is an interesting relationship between the degree of a node and the correlation of its weighted clustering coefficient with its trustfulness.

6.3. Comparison Between Different Definitions

In this subsection, we compare the three definitions of weighted clustering coefficient with the classic definition of unweighted clustering coefficient. We show that the new

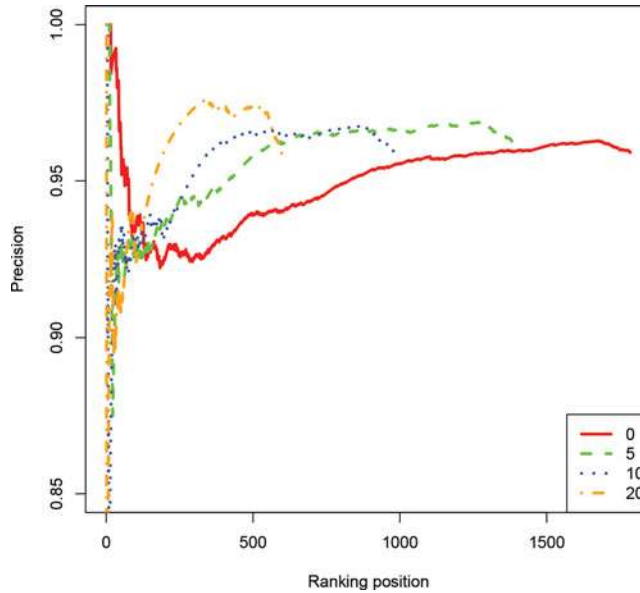


Figure 4 Precision of the sample algorithm when we restrict only to nodes with unweighted degree bigger than 0, 5, 10 and 20 for $W C_v^{\text{random}}$.

definition is always comparable with the other two and at various points of the ranking, it performs significantly better.

In Figure 5 we show the ranking obtained using the four definitions. For the classic unweighted definition we compute the clustering coefficient of each node exactly. For the three weighted clustering coefficient definitions we approximate the clustering coefficient using 3200 samples per node.

It is possible to note that the two previous definitions of clustering coefficient have very similar performances and performances very similar to the classic unweighted definition, whereas the ranking obtained by our new definition has higher precision for the first positions in the ranking and then has performances comparable with the rankings obtained using the other definitions.

Finally, we compare the performances of our new definition with the performances of the definition given by Barrat et al. when we restrict computation to nodes with unweighted degree above a specific threshold. We think that this case is of particular interest because we showed in the previous subsection that there is an interesting relationship between the degree of a node and the correlation between its weighted clustering coefficient and its trustfulness.

In Figure 6, we present the comparison between the two definitions when we restrict our attention to nodes of degree larger than 0, 5 and 20.

Also in this case, we observe that the two definitions have very similar performances.

6.4. Scalability of Our Algorithm

In this final subsection, we analyze the scalability of our new algorithm. In order to do it, we analyze five public datasets available in the SNAP repository: Orkut, Patents, Pokec,

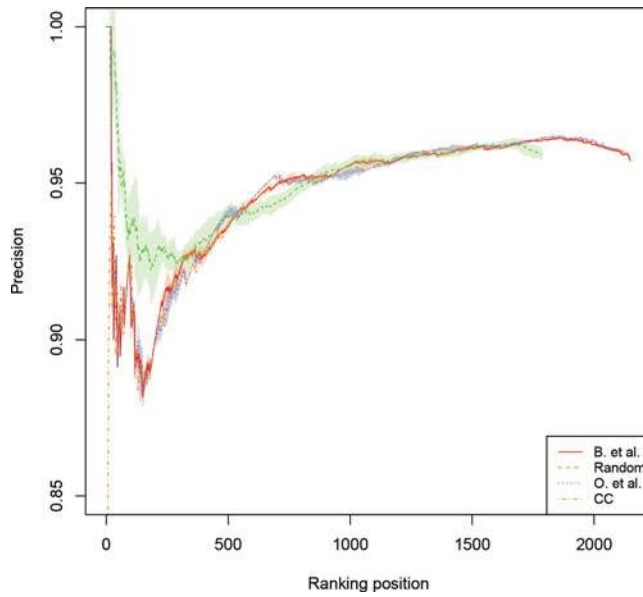


Figure 5 Comparison between the precision rankings obtained by classic definition of clustering coefficient (CC), the definition by Barrat et al. (B. et al.), the definition of Onnela et al. (O. et al.), and our new definition (Random).

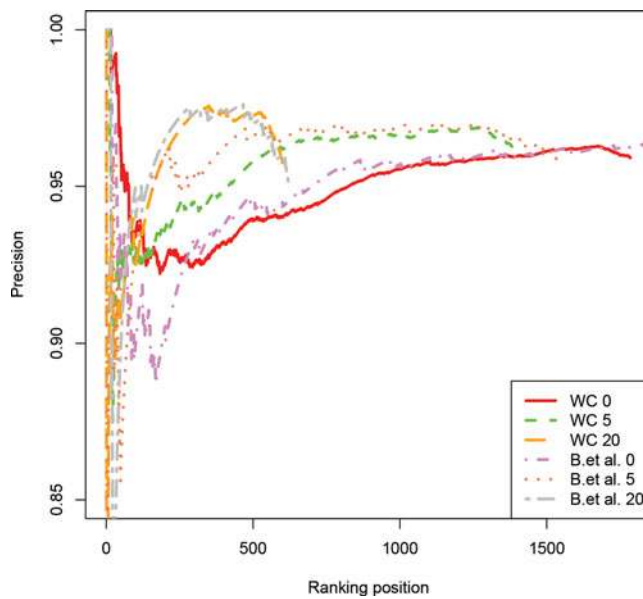


Figure 6 Comparison between the precision rankings obtained by the definition by Barrat et al. (B. et al.) and our new definition (WC) when we restrict computation to node of degree larger than 0, 5 and 20.

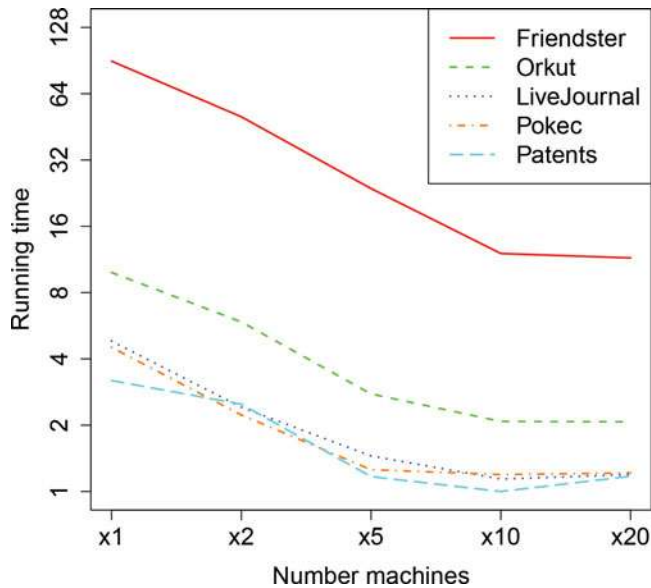


Figure 7 Running time of the algorithm on different networks with different resources. Note that all the numbers shown in the plots are relative. In particular, on the y axis we present the running times as the relative running time in comparison with the fastest run of algorithm on the smallest graph, and on the x axis we present the number as a multiplicative factor of the minimum number of machines used $x2$, $x5$, $x10$, and $x20$.

LiveJournal and Friendster. The five graphs have increasing numbers of edges (each dataset has roughly twice as many edges as the previous one). The input dataset are unweighted, so before we run our algorithm, we assign a random weight between 0 and 1 to every edge.

In this experiment we are interested in analyzing the running time of the algorithm when we increase the number of machines available and when we increase the size of the network analyzed.

In Figure 7, we present the running time of the algorithm on different networks with different resources. Note that all the numbers shown in the plots are relative. In particular, on the y axis we present the running times as the relative running time in comparison with the fastest run of algorithm on the smallest graph, and on the x axis we present the number as a multiplicative factor of the minimum number of machines used.

It is interesting to note that our algorithm is able to leverage on parallelization to speed-up computation on very large graphs. Note, for example, that by increasing the number of machines by a factor of 10, it is possible to reduce the running time on the Friendster graph of roughly a factor of 10. It is also worth noticing that we do not obtain much gain by increasing the number of machines by a factor of 20; this is probably due to a trade-off between computational power and cost of communication within machines.

7. CONCLUSIONS

In this work we present sampling techniques to obtain efficient estimators for several measures of weighted clustering coefficient together with their MapReduce implementation. We also propose a novel graph-theoretic notion of clustering coefficient in weighted networks defined as the expected unweighted clustering coefficient on a family of random

graphs. Moreover, we show on an application related to web spam detection that the notions of weighted clustering coefficient compare with the standard notion of unweighted clustering coefficient as a machine learning feature to assess the quality of nodes in a social network. Given the importance of weighted networks to model the strength of the interaction between nodes in a graph, we hope that our work will prompt more study on the relevance of weighted graph mining features to characterize the inner structure of social networks.

ACKNOWLEDGMENTS

We thank Corinna Cortes for suggesting the problem.

FUNDING

Part of this work was done while visiting Google Research, NY. Partly supported by Google Focused award “Web Algorithmics for Large-scale Data Analysis” and ERC StG PAAI “Practical Approximation Algorithms.”

REFERENCES

- [1] L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan. “Group Formation in Large Social Networks: Membership, Growth, and Evolution.” Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadelphia, PA, USA, August 20–23, 2006. ACM 2006.
- [2] A. Barrat, M. Barthélemy, R. Pastor-Satorras, and A. Vespignani. “The Architecture of Complex Weighted Networks.” In *Proceedings of the National Academy of Sciences of the United States of America*, PNAS 2004 101 (11) 3747–3752.
- [3] L. Becchetti, P. Boldi, C. Castillo, and A. Gionis. “Efficient Semi-Streaming Algorithms for Local Triangle Counting in Massive Graphs.” Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Las Vegas, Nevada, USA, August 24–27, 2008. ACM 2008.
- [4] B. Bollobas. “Mathematical Results on Scale-Free Random Graphs.” In *Handbook of Graphs and Networks: From the Genome to the Internet*, edited by S. Bornholdt and H. G. Schuster. Wiley-VCH Verlag GmbH & Co. KGaA, Weinheim, FRG.
- [5] F. Bonchi, F. Gullo, A. Kaltenbrunner, and Y. Volkovich. “Core Decomposition of Uncertain Graphs.” In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’14, pp. 1316–1325. New York, NY, USA: ACM, 2014.
- [6] C. Budak, D. Agrawal, and A. El Abbadi. “Structural Trend Analysis for Online Social Networks.” *PVLDB* 4(10): 646–656 (2011).
- [7] L. Buriol, G. Frahling, S. Leonardi, A. Marchetti-Spaccamela, and C. Sohler. “Counting Triangles in Data Streams.” Proceedings of the Twenty-Fifth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 26–28, 2006, Chicago, Illinois, USA. ACM 2006.
- [8] C. Castillo, D. Donato, L. Becchetti, P. Boldi, S. Leonardi, M. Santini, and S. Vigna. “A Reference Collection for Web Spam.” *SIGIR Forum* 40(2): 11–24 (2006).
- [9] J. Dean and S. Ghemawat. “MapReduce: Simplified Data Processing on Large Clusters.” *Communications of the ACM* 51(1): 107–113 (2008).
- [10] G. Fagiolo. “Clustering in Complex Directed Networks.” *Phys. Rev. E* 76 026107 – Published 16 August 2007.

- [11] S. J. Hardiman and L. Katzir. "Estimating Clustering Coefficients and Size of Social Networks via Random Walk." 22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13–17, 2013. International World Wide Web Conferences Steering Committee / ACM 2013, pp. 539–550.
- [12] P. Hintsanen and H. Toivonen. "Finding Reliable Subgraphs from Large Probabilistic Graphs." *Data Min. Knowl. Discov.* 17(1): 3–23 (2008).
- [13] M. Jha, C. Seshadhri, and A. Pinar. "A Space Efficient Streaming Algorithm for Triangle Counting Using the Birthday Paradox." The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11–14, 2013. ACM 2013, pp. 589–597.
- [14] G. Kalna and D. J. Higham. "Clustering Coefficients for Weighted Networks." In *Symposium on Network Analysis in Natural Sciences and Engineering*. Bristol, UK, April 5–6, 2006.
- [15] H. Kwak, C. Lee, H. Park, and S. Moon. "What is Twitter, a Social Network or a News Media?" Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26–30, 2010. ACM 2010, pp. 591–600.
- [16] M. Latapy. "Main-Memory Triangle Computations for Very Large (Sparse(Power-Law)) Graphs." *Theoretical Computer Science* 407(1–3): 458–473 (2008).
- [17] J. Leskovec and E. Horvitz. "Planetary-Scale Views on a Large Instant-Messaging Network." Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21–25, 2008. ACM 2008, pp. 915–924.
- [18] J. Leskovec, J. Kleinberg, and C. Faloutsos. "Graphs over Time: Densification Laws, Shrinking Diameters and Possible Explanations." Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Chicago, Illinois, USA, August 21–24, 2005. ACM 2005, pp. 177–187.
- [19] E. Liberty. "Simple and Deterministic Matrix Sketches." The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11–14, 2013. ACM 2013, pp. 581–588.
- [20] L. Lopez-Fernandez, G. Robles, and J. M. Gonzalez-Barahona. "Applying Social Network Analysis to the Information in cvs Repositories." *International Workshop on Mining Software Repositories (MSR 2004)*, 2004 p. 101–105. Edinburgh, United Kingdom, May 23–28, 2004.
- [21] M. E. J. Newman. "Analysis of Weighted Networks." *Phys. Rev. E* 70 (Nov 2004), 056131.
- [22] M. E. J. Newman, D. J. Watts, and S. H. Strogatz. "Random Graph Models of Social Networks." *Proc. Natl. Acad. Sci. USA* 99 (2002), 2566–2572.
- [23] J.-P. Onnela, J. Saramäki, J. Kertész, and K. Kaski. "Intensity and Coherence of Motifs in Weighted Complex Networks." *Physical Review E*. 71, 065103(R) – Published 13 June 2005.
- [24] T. Opsahl and P. Panzarasa. "Clustering in Weighted Networks." *Social Networks*. 32:245–251, 2010.
- [25] R. Pagh and C. E. Tsourakakis. "Colorful Triangle Counting and a MapReduce Implementation." *Inf. Process. Lett.* 112(7): 277–281 (2012).
- [26] J. Saramäki, M. Kivelä, J.-P. Onnela, K. Kaski, and J. Kertesz. "Generalizations of the Clustering Coefficient to Weighted Complex Networks." *Physical Review E* 75, 027105 – Published 23 February 2007.
- [27] T. Schank and D. Wagner. "Approximating Clustering Coefficient and Transitivity." *J. Graph Algorithms Appl.* 9(2): 265–275 (2005).
- [28] T. Schank and D. Wagner. "Finding, Counting and Listing all Triangles in Large Graphs, an Experimental Study." Experimental and Efficient Algorithms, 4th International Workshop, WEA 2005, Santorini Island, Greece, May 10–13, 2005, Proceedings. Lecture Notes in Computer Science 3503, Springer 2005, pp. 606–609.
- [29] S. Suri and S. Vassilvitskii. "Counting Triangles and the Curse of the Last Reducer." Proceedings of the 20th International Conference on World Wide Web, WWW 2011, Hyderabad, India, March 28 – April 1, 2011. ACM 2011, pp. 607–614.

- [30] L. Takac and M. Zabovsky. "Data Analysis in Public Social Networks." In *International Scientific Conference & International Workshop Present Day Trends of Innovations*, Łomża, Poland, May 28–29, 2012.
- [31] C. E. Tsourakakis, U. Kang, G. L. Miller, and C. Faloutsos. "Doulion: Counting Triangles in Massive Graphs with a Coin." Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Paris, France, June 28 – July 1, 2009. ACM 2009, pp. 837–846.
- [32] C. E. Tsourakakis, M. N. Kolountzakis, and G. L. Miller. "Triangle Sparsifiers." *J. Graph Algorithms Appl.* 15(6): 703–726 (2011).
- [33] D. J. Watts and S. H. Strogatz. "Collective Dynamics of Small-World Networks." *Nature* 393, 440–442 (1998).
- [34] J. Yang and J. Leskovec. "Defining and Evaluating Network Communities Based on Ground-Truth." *Knowl. Inf. Syst.* 42(1): 181–213 (2015).
- [35] Z. Yang, C. Wilson, X. Wang, T. Gao, B. Y. Zhao, and Y. Dai. "Uncovering Social Network Sybils in the Wild." *TKDD* 8(1): 2:1–2:29 (2014).
- [36] B. Zhang, S. Horvath, et al. "A General Framework for Weighted Gene Co-Expression Network Analysis." *Stat Appl Genet Mol Biol.* 2005;4:Article17.
- [37] Y. Zhang, Z. Zhang, J. Guan, and S. Zhou. "Analytic Solution to Clustering Coefficients on Weighted Networks." *Journal of Statistical Mechanics: Theory and Experiment*, 2010.