
Preserving Pairwise Relationships in Subgraphs

Andrea Vattani

UC San Diego, 9500 Gilman Drive, La Jolla, CA 92093

AVATTANI@CS.UCS.D.EDU

Deepayan Chakrabarti

Yahoo! Research, 701 1st Avenue, Sunnyvale, CA 94089

DEEPAY@YAHOO-INC.COM

Maxim Gurevich

Yahoo! Research, 701 1st Avenue, Sunnyvale, CA 94089

MAXIMG@YAHOO-INC.COM

Abstract

Choosing a subgraph that can concisely represent a large real-world graph is useful in many scenarios. The usual strategy employed is to sample nodes so that the *induced* subgraph matches the original graph's degree distribution, clustering coefficient, etc., but no attempt is made to preserve fine-grained relationships between nodes, which are vital for applications like clustering, classification, and ranking. In this work, we model such relationships via the notion of Personalized PageRank Value (PPV). We show that induced subgraphs output by current sampling methods do not preserve PPVs, and propose algorithms for creating PPV-preserving subgraphs on *any* given subset of graph nodes. Experiments on three large real-world graphs show that the subgraphs created by our method improve upon induced subgraphs in terms of preserving PPVs, clustering accuracy, and maintaining basic graph properties.

1. Introduction

Consider the problem of clustering a large graph or finding communities in a social network. There are many available methods, and their performance varies considerably across graphs, so picking the right method for a given graph is important. Not only that, but each method has its own parameter settings

that need to be tuned (e.g., the conductance threshold for (Andersen et al., 2008), or the edge-weight transform function for (van Dongen, 2000)). On a small-sized dataset, the obvious solution is to search the space of parameters iteratively until a good parameter setting is found. However, each iteration requires running the clustering algorithm on the entire graph and observing the quality of the solution. Moreover, this tuning must be performed for every available clustering method to pick the best one. Clearly, this does not scale well to large graphs.

One solution is to find smaller graphs $\{S_1, S_2, \dots\}$ that mimic the full graph G , and perform tuning and model selection on these S_i . For instance, if we knew the generative model of G , then we could generate S_i from this model. However, real-world graphs rarely match any generative model in all respects, and even then, finding the model parameters might be too expensive.

Alternatively, we could build small graph *samples* from G itself. There has been much recent work on picking “representative” samples of a graph that preserve degree distributions, clustering coefficients, and so on (Leskovec & Faloutsos, 2006; Maiya & Berger-Wolf, 2010). However, these measures might not be relevant for the application (say, clustering), and the measures that are critical to the application (say, pairwise similarities between nodes, conductance between node subsets, etc.) might not be preserved in the sample. Thus, the challenges are twofold: (a) identify a set of measures that are important for a wide range of applications, and (b) design a sampling method that preserves these particular measures, in addition to usual ones like degree distributions.

A COMMON MEASURE. Random walk based algorithms have been used in many common graph appli-

Appearing in *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, WA, USA, 2011. Copyright 2011 by the author(s)/owner(s).

cations: propagating labels in classification (Szummer & Jaakkola, 2001), finding spam webpages (Gyöngyi et al., 2004), clustering (Andersen et al., 2008), search in peer-to-peer networks (Gkantsidis et al., 2006), collaborative filtering applications (Gori & Pucci, 2007), ranking (Haveliwala, 2003), etc. In addition, random walks have been shown to be strictly related to other fundamental graph quantities, such as normalized cuts (Meila & Shi, 2000) and the graph spectrum (Orponen et al., 2008). This motivates our goal of creating graph samples that “preserve random walks”, that is preserve probabilities p_{ij} of a walk starting from node i hitting node j , for all i, j in the sample graph. These are in fact the Personalized Pagerank Values or PPVs (ignoring the restart probability), and they reflect the pairwise similarities between nodes. Indeed, in all the applications above, random walks are used to compute such similarities (possibly implicitly), which are in turn used for clustering, classification, ranking, etc. Thus, our goal is to build a graph sample where the PPVs between all pairs of nodes in the sample are preserved. Such samples are useful not only for random-walk applications, but also for applications optimizing seemingly unrelated measures — as we show empirically in Section 6.

SAMPLING METHOD. Graph sampling typically consists of two parts: (a) picking a subset S of nodes from the full graph G , and (b) selecting the edges between nodes in S . Current sampling algorithms focus almost exclusively on part (a); for part (b), they simply retain those edges in G both of whose endpoints fall in S (call this the “induced” subgraph). Unfortunately, induced subgraphs often fail to preserve PPVs. For instance, two nodes in S that had many common neighbors in G but were not directly linked may become unreachable from each other in the induced subgraph, even if intuitively they are “close” in G and ought to be “close” in the sample as well. Since PPV between these nodes is likely to be high, they would remain “close” in a PPV-preserving subgraph. This motivates us to focus on part (b), and in particular on *creating* a subgraph that preserves PPVs between all pairs of nodes in S , whose edges might, or not, have existed in G . A major advantage of our method is that it applies to *any* subset S , so we can leverage all previous algorithms developed for part (a).

To illustrate the above we show, in Figures 1-3, a graph of about 2,000 nodes, the induced subgraph on a 200-node sample, and the subgraph created by our algorithm on the same sample and having the same average node degree as the full graph. Observe that the induced subgraph only preserves edges between imme-

diately neighbors, leading to many disconnected components. The subgraph generated by our algorithm better preserves graph structure by creating “bridging” edges between components that are most tightly connected (according to PPV) in the original graph.

RESULTS. Our contributions are the following:

- (1) We show that induced subgraphs do not preserve PPVs (Section 4), and that even the simpler problem of preserving global PageRank in induced subgraphs is difficult.
- (2) We show that, while it might not be possible in general to preserve PPVs for a given subset, it is possible to *exactly* match PPVs by adding just one extra “sink” node to this subset. We present algorithms implementing this approach on an *arbitrary* subset of nodes of a directed graph (Section 5.1). Additionally, we show that adding an extra “source” node allows to exactly maintain global PageRank values as well (Section 5.3).
- (3) Perfectly preserving PPVs may require weighted edges between all pairs of nodes in the subset. To reduce the storage requirements, we present rounding methods that remove a large fraction of these edges while still providing provable bounds on the distortion of PPVs (Section 5.2).
- (4) Finally, we verify via experiments on three large real-world networks that subgraphs generated by our algorithms outperform existing methods in terms of preserving PPVs, improving clustering accuracy, and preserving basic graph properties. (Section 6).

Due to space constraints, complete proofs of our claims will appear in the full version of the paper.

2. Related Work

Sampling from large graphs while preserving graph properties was introduced by Leskovec and Faloutsos (2006). They proposed several criteria for “good” graph samples and empirically evaluated several node sampling algorithms with respect to these criteria. They observed that different algorithms are better at preserving some graph properties but not others, and that there is no clear winner among the evaluated algorithms. Another line of work considers sampling subgraphs to optimize for one specific goal, such as the preservation of community structures (Maiya & Berger-Wolf, 2010), or the ability to visualize the graph (Jia et al., 2008). Hubler et al. (2008) proposed a generic Metropolis algorithm for sampling nodes so as to minimize distortion of a given graph property. However, the number of steps until convergence is not known in advance (and may be large), and each step requires the evaluation of the distortion function over

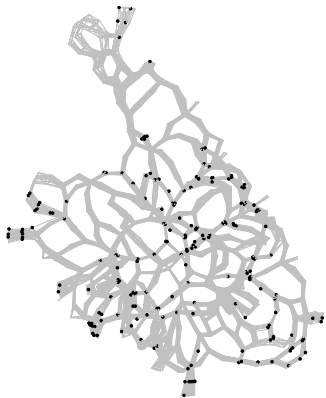


Figure 1. The original graph. The 200 sampled nodes are black.

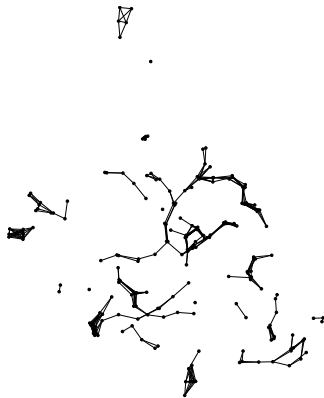


Figure 2. The induced subgraph on the sample.

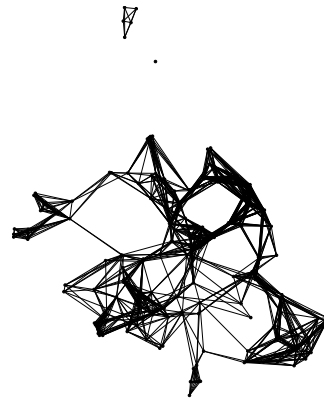


Figure 3. The subgraph produced by our algorithm (UNWEIGHTED).

numerous subgraphs, which may be costly. Our work deviates from all of these in that we are not concerned with sampling of nodes, but rather, given a subset of nodes, we consider the general problem of maintaining multi-scale graph structure by preserving a distance metric (personalized PageRank) between all pairs of sampled nodes.

Instead of sampling, Sala et al. (2010) proposed generating synthetic graphs that have similar properties to the original graph. While having advantages like privacy preservation, this approach tries to solve the harder problem of generating a small representative graph using only few parameters, and thus needs to assume a model of the original graph. Our sampling based approach uses more information about the original graph and thus is applicable to graphs that do not accurately follow any known model.

Another related area of research is graph sparsification. Edge sparsification only decreases the number of edges, leaving the set of nodes unchanged. Thus it is not applicable when the goal is to obtain a significantly smaller representative graph. In vertex sparsification, the number of vertices is reduced. Moitra (2009) proposed an algorithm that, given a graph and a subset of nodes, constructs a new graph on these nodes that approximately preserves minimum (unnormalized) cuts. In many applications such as clustering or classification, PPV-related measures like normalized cuts and conductance are believed to be much more effective than unnormalized cuts (Kannan et al., 2004). We thus focus on preserving personalized PageRank which better reflects the relative importances of nodes.

3. Preliminaries

Consider a directed graph $G = (V, E)$ of n nodes, its adjacency matrix A and diagonal out-degree matrix

D . The PageRank of G is defined as the vector \mathbf{p} that solves the following system:

$$\mathbf{p} = \alpha \mathbf{r} + (1 - \alpha) A^t D^{-1} \mathbf{p}, \quad (1)$$

where $0 < \alpha < 1$ is the restart probability (typically, $\alpha \approx 0.15$) and \mathbf{r} is the personalization vector. Uniform $\mathbf{r}(i) = \frac{1}{n}$ for all $i \in V$ results in \mathbf{p} being the *global* PageRank. If \mathbf{r} contains 1 in the i -th coordinate and is 0 elsewhere, the solution of Eq. 1 gives the i -personalized PageRank \mathbf{p}_i (Haveliwala, 2003). From such \mathbf{p}_i 's and the linearity of Eq. (1), the solution \mathbf{p} for arbitrary \mathbf{r} can be easily derived. A $\mathbf{p}_i(j)$ can also be thought of as the probability of an α -discounted random walk from node i ending in node j :

$$\mathbf{p}_i(j) = \sum_{t=0}^{\infty} \alpha (1 - \alpha)^t \mathbf{p}_i^{[t]}(j), \quad (2)$$

with $\mathbf{p}_i^{[t]}(j)$ being defined as

$$\mathbf{p}_i^{[t]}(j) = \sum_{k_1=i, k_2, \dots, k_{t+1}=j} \prod_{l=1}^t \frac{1}{d^+(k_l)},$$

where $d^+(k_\ell)$ is the out-degree of node k_ℓ , and the summation is over all paths of length t from i to j . An equivalent formulation is given by the decomposition theorem (Jeh & Widom, 2003):

$$\mathbf{p}_i = \alpha \mathbf{e}_i + (1 - \alpha) \sum_{j:(i,j) \in E} \frac{\mathbf{p}_j}{d^+(i)}.$$

The decomposition theorem suggests that the ‘‘importance’’ that a node i gives to node $j \neq i$ is proportional to the average importance that i 's out-neighbors give to j . All of these can be easily extended to the case of weighted graphs by replacing inverse out-degrees by normalized weights, i.e., $1/d^+(i) = w(i, j) / \sum_j w(i, j)$, where $w(i, j)$ is the weight of edge (i, j) .

4. Induced subgraphs

In this section we examine the problem of finding good representative samples using induced subgraphs. The quality of a subgraph is measured by how well the PageRanks of the nodes in the subgraph are preserved.

Consider an undirected graph $G = (V, E)$, where the degree of node i is denoted by $\deg_G(i)$, and restart probability $0 < \alpha < 1$. Note that α induces a global PageRank vector \mathbf{p}^G ($\mathbf{p}^G(i)$ is the PageRank of i). Suppose we are given a sample of size $0 < k < |V|$ and are interested in finding a subset $S \subset V$ of nodes of size k such that the PageRanks in the subgraph $G[S]$ induced by S are as close as possible to \mathbf{p}^G for all nodes in S . By looking at $\mathbf{p}^{G[S]}$ and \mathbf{p}^G as vectors, we want to minimize some notion of distance between them. Two remarks are necessary: first, \mathbf{p}^G is $|V|$ -dimensional, so we only consider its restriction to nodes in S ; second, since $\sum_{i \in V} \mathbf{p}^G(i) = 1$ and $\sum_{i \in S} \mathbf{p}^{G[S]}(i) = 1$, we rescale \mathbf{p}^G ; therefore we consider the $|S|$ -dimensional scaled vector $\tilde{\mathbf{p}}^G$, where $\tilde{\mathbf{p}}^G(i) := \mathbf{p}^G(i) / \sum_{j \in S} \mathbf{p}^G(j)$ for all $i \in S$. Now our problem is finding an induced subgraph S of G of size k such that $d(\tilde{\mathbf{p}}^G, \mathbf{p}^{G[S]})$ is minimized, where d is some distance function (e.g. a norm).

We show that this problem is NP-hard for any non-trivial distance function d : we only require that $d(\mathbf{a}, \mathbf{a}) = 0$ and that $d(\mathbf{a}, \mathbf{b}) > 0$ if $\mathbf{a} \neq \mathbf{b}$.

Theorem 1. *For any $0 < \alpha < 1$, it is NP-hard to find a subset $S \subset V$, $|S| = k$, that minimizes $d(\tilde{\mathbf{p}}^G, \mathbf{p}^{G[S]})$ for any (non-trivial) distance function d . The NP-hardness holds even when restricted to cubic graphs.*

Next, we show that there exists a family of graphs where a random subset S is likely to contain a node that has higher PageRank relative to many other nodes in S when computed on the induced subgraph, but lower relative PageRank when computed over the entire graph G . Clearly, such induced subgraphs would fail to preserve PPVs as well. The proof of the theorem uses concentration bounds on the degrees of the nodes in a (Erdős–Rényi) random graph $G(n, p)$.

Theorem 2. *Let $0 < p < 1$ and $G = (V, E)$ be sampled from $G(n, p)$. Let S be a random subset of V with $|S| = o(n/\log n)$. Then, with probability at least $1/2 - o(1)$, there exists an $S' \subset S$ of size at least $\Omega(|S|)$, such that a node u^* of maximum degree in $G[S]$ has degree smaller than every node in S' in G . Therefore, for small enough α , every node in S' has PageRank higher than u^* in G , and lower than u^* in $G[S]$.*

5. Preserving PPV

The previous section showed that preserving even *global* graph properties such as PageRank on *induced* subgraphs is difficult. This is expected, since induced subgraphs only preserve links between immediate neighbors, ignoring multi-hop connections. In this section we relax the problem definition by allowing the creation of edges that did not exist in the original graph, and show that this yields a much better approximation of node-personalized PageRank (PPV).

For the rest of this section we assume that all non-existent edges in G have zero weight, and that outgoing degrees are normalized, i.e., $\sum_{j \in V} w_G(i, j) = 1$ for every $i \in V$. This is without loss of generality since PPVs depend only on the normalized weights. We first show that PPVs can be preserved *exactly* using a weighted graph construction that adds a special “sink” node, and present an algorithm to compute edge weights. However, this graph can contain up to $|S|^2$ edges between all pairs of nodes in S . To alleviate this, we present two sparsification algorithms, one of which yields a sparse weighted graph, and the other a sparse unweighted graph. We bound the error in PPVs introduced by this sparsification. Finally, we show that with the addition of another special “source” node, it is possible to preserve global PageRanks in addition to all pairwise PPVs. We stress that our subgraph generation method can be applied to *any* subset of nodes, so popular node-sampling methods in the literature (Leskovec & Faloutsos, 2006) can be used as a first step to pick a node sample on which our algorithms are applied.

5.1. Constructing a PPV-preserving subgraph

Suppose we are given a weighted directed graph $G = (V, w_G)$, where the edge weight function $w_G : V \times V \rightarrow \mathcal{R}^+$ defines the edge set, and a subset $S \subset V$ of nodes. We consider the problem of constructing a new weighted graph $H = (S, w_H)$ only on the subset S that preserves PPVs between nodes in S – that is, all values $\mathbf{p}_i^H(j) = \mathbf{p}_i^G(j)$ for $i, j \in S$.

Observation 3. *There exists a graph $G = (V, E)$ and a subset $S \subset V$ such that no weighted graph $H = (S, w_H)$ (with non-negative weights) preserves all PPVs between nodes in S .*

This can be shown via a linear program that reaches an optimum with zero error if it is possible to preserve all PPVs exactly, and a positive error otherwise. Examples of subsets with non-zero error can then easily be found.

Despite the negative result of Obs. 3, adding an extra “sink” node, having only incoming edges from S , allows to preserve PPVs of nodes in H . The main result of this section is the following.

Theorem 4. *Let $G = (V, w_G)$ be a directed weighted graphs. For any $S \subset V$, there exists a weighted directed graph H on $S \cup \{\text{SINK}\}$ where PPVs between nodes in S are preserved: for all $i, j \in S$, $\mathbf{p}_i^H(j) = \mathbf{p}_i^G(j)$. Moreover, there exists a polynomial time algorithm that computes such H given G .*

The idea of the algorithm, presented in Table. 1, is to start with the original graph G and remove nodes in $V \setminus S$ one by one. For each node z being removed, the algorithm identifies the set of its neighbors, $N(z)$, removes all edges incident to z , and creates edges between all pairs of its neighbors, and from each neighbor to the sink. Specifically, for every $j, k \in N(z)$, directed edges (j, k) and (k, j) are created (parallel edges are merged). The weight of edge (j, k) captures the probability of a random walk starting at j to reach k through z . These new edges, however, do not capture the probability of the walk stopping at z , so an edge to the sink is created from each $j \in N(z)$.

Algorithm 1 NODEREMOVAL

Input: $G = (V, w_G), S \subseteq V$.

Output: $H = (S \cup \{\text{SINK}\}, w_H)$ s.t. $\mathbf{p}_i^H(j) = \mathbf{p}_i^G(j)$, $i, j \in S$.

Initialize H to be equal to G

Add a node SINK in H : $w_H(i, \text{SINK}) = 0$ for all i , and $w_H(\text{SINK}, \text{SINK}) = 1$

for each $z \in V \setminus S$ **do**

 Remove z from H (and all edges incident to it)

for each node $x \neq z$ such that $w_G(x, z) > 0$ **do**

for each node $y \neq z$ such that $w_G(z, y) > 0$ **do**

 Add a new edge $e^*(x, y)$ in H , and set
 $w_H(e^*(x, y)) += (1 - \alpha)w_G(x, z) \cdot w_G(z, y)$
 $\cdot \sum_{t=0}^{\infty} [(1 - \alpha)w_G(z, z)]^t$

end for

 Add a new edge $e^*(x, \text{SINK})$ in H , and set

$w_H(e^*(x, \text{SINK})) +=$
 $w_G(x, z) - \sum_{y \neq z: w_G(z, y) > 0} w_H(e^*(x, y))$

end for

 Set G (and its weights) to be H

end for

Thm. 4 follows from the next two lemmas. Lemma 5 follows from simple algebraic manipulations.

Lemma 5. *In each iteration of NODEREMOVAL, (a) all added edges have non-negative weights, and (b) the (weighted) out-degree of each node sums to 1.*

Lemma 6. *An iteration of NODEREMOVAL preserves PPVs between all but the removed node.*

Proof Sketch. The idea of the proof is to provide a mapping between the contributions to PPVs by walks in G and the contributions to the same PPVs by walks in H . The proof shows that a single walk through a newly added edge in H “summarizes” contributions of countably infinitely many walks that pass through the removed node z in G .

PRESERVING PPVs ON SMALL SUBSETS. While the NODEREMOVAL algorithm can be applied to all graphs, it can be costly for small subsets, for which many nodes need to be removed. For this case, we present an alternative algorithm below.

Recall that $\mathbf{p}_i^H(j)$ is the desired PPV value from node i to j , where $i, j \in S$, and we want this value to be equal to $\mathbf{p}_i^G(j)$. Let P^H be the matrix where $\mathbf{p}_i^H(\cdot)$ is row i . Let W_H be the edge weight matrix in H , i.e., $w_H(i, \cdot)$ is row i in W_H . P^H is known, and we want to find W_H (or equivalently, its transpose W_H^t).

Lemma 7. *We have the following relation:*

$$W_H^t = \frac{1}{1 - \alpha} (I - \alpha(P^H)^{-1})$$

This follows from manipulations of equation (1). The existence of a solution is guaranteed by Thm. 4. This also implies the uniqueness of the solution.

The initial computation of pairwise PPVs takes $O(E + |S|^2)$ time (using (Sarlós et al., 2006)). This lets us compute W_H via inversion of a matrix of size $(|S| + 1)$ at a complexity of $O(|S|^{2.376})$. For small subsets, this can be significantly cheaper than NODEREMOVAL.

5.2. Rounding the Weighted Matrix

The weighted graph H obtained via NODEREMOVAL or Lemma 7 could be a dense graph with up to $(k + 1)^2$ edges. We may want a much sparser subgraph to reduce storage requirements. Next, we present two algorithms for sparsifying the graph, one generating weighted, and the other unweighted graphs. Still, both maintain weighted edges to the sink (which can be thought of as storing one extra number per node in S); as we will later show in Section 6, these edges are important for maintaining PPVs.

Both algorithms are simple instances of rounding: we set a threshold τ and for all edges (u, v) such that $w_H(u, v) < \tau$, we increase $w_H(u, \text{SINK})$ by $w_H(u, v)$ and then set $w_H(u, v) = 0$. Edges with zero weight can obviously be removed, thus sparsifying the graph.

Note that this leaves the weights of the remaining edges unchanged, and all node out-degrees still sum to one. We call this algorithm **WEIGHTED**.

The second algorithm (called **UNWEIGHTED**) takes the output of **WEIGHTED** and makes all edges unweighted (except for the edges to the sink). Thus, **WEIGHTED** and **UNWEIGHTED** give sparse weighted and unweighted graphs respectively (apart from the weights on edges to sink).

Next, we show that the PPV matrix (as defined before Lemma 7) for a graph rounded by **WEIGHTED** is close to the PPV matrix for the unrounded graph (where the original PPVs are preserved exactly). Let H denote the subgraph obtained via **NODEREMOVAL**, and R the subgraph after applying **WEIGHTED** on H . Also, for any node j in S , let $P^H(j)$ and $P^R(j)$ be the j -th columns of P^H and P^R respectively.

Theorem 8. *For all $j \in S$,*

$$\frac{\|P^H(j) - P^R(j)\|_1}{|S| + 1} \leq 2\tau \frac{1 - \alpha}{\alpha}.$$

The proof follows from a coupling between the Markov chains representing the full PPV matrix and the rounded one. Thus, the rounded graph created by **WEIGHTED** provides a good approximation to the PPVs of the original graph, provided τ is not too high.

We remark here that in the unlikely boundary case of a node with extremely high out-degree, each of its outgoing edges might get a small weight and be rounded, leading to a seeming paradox of a high-degree node being rounded to a node of zero degree. However, this implies a large $|S|$ and so the bound will be weak unless τ is correspondingly small as well. Also, the edges that reflect the importance of a node are its incoming edges, and these will still be preserved. In our experiments, we pick τ so as to ensure that the rounded subgraph has as many edges as the induced subgraph; this caused no extreme roundings on our graphs. However, if tighter bounds are desired, a smaller τ must be chosen (and hence, more edges retained).

5.3. Preserving Global PageRanks

In this section we show that it is possible to preserve the global PageRank of each node in S by adding to H (computed in Section 5, without rounding) a single “source” node and edges from it to all nodes in S , or equivalently, storing one extra number per node in S .

The following corollary shows that as long as the nodes removed by **NODEREMOVAL** have same “weight” in the personalization vector \mathbf{r} , which is the case for global

PageRank, a single value per each node in S can summarize the contributions of PPVs from $V \setminus S$ to the node in S . Let \mathbf{v} be an $|S|$ -dimensional vector containing these contributions: $\mathbf{v}(i) = \sum_{j \in V \setminus S} \mathbf{p}_j^G(i)$, $i \in S$. From the linearity of PageRank (Equation (1)) we get:

Corollary 9. *Suppose we are given H , constructed from G by algorithm **NODEREMOVAL**, and the vector \mathbf{v} . Then, for every personalization vector \mathbf{r} such that $\mathbf{r}(j) = \mathbf{r}(j')$ for $j, j' \in V \setminus S$, \mathbf{r} -personalized PageRank $\mathbf{p}_r^G(i)$ can be computed for all $i \in S$.*

The information contained in \mathbf{v} can be naturally integrated into the subgraph by adding a single node we call **SOURCE**. Let H' be the graph obtained from H by adding a node **SOURCE**, and weighted edges from **SOURCE** to other nodes in H (including **SINK**). Then, standard techniques on H' can be used to compute the global PageRank.

To compute the weights of the new edges from **SOURCE**, we use an auxiliary graph G' capturing the partial PPV contributions of random walks in G that start outside S and end in S , without ever hitting S in the middle. Intuitively, this accounts for the PPV “mass” that flows into S ; the distribution of this mass within S and the loss of mass from S to $V \setminus S$ is already handled by our subgraph on $S \cup \{\mathbf{SINK}\}$. In particular, G' is obtained from G by (a) removing all edges from S to $V \setminus S$, (b) collapsing all the nodes in S into a single node, (c) summing up weights of parallel edges, and (d) adding a self-loop of weight 1 to the collapsed node. Then, for $\ell \in S$, edge weights are set to:

$$w_{H'}(\mathbf{SOURCE}, \ell) = \frac{1}{|V \setminus S|} \frac{1}{\alpha} \sum_{j, k \in V \setminus S} \mathbf{p}_j^{G'}(k) \cdot w_G(k, \ell).$$

It can be shown that the weighted out-degree of **SOURCE** sums to one or less. To normalize it we add an edge to **SINK**: $w_{H'}(\mathbf{SOURCE}, \mathbf{SINK}) = 1 - \sum_{i \in S} w_{H'}(\mathbf{SOURCE}, i)$. Now, in a manner similar to Thm. 4, we may prove:

Theorem 10. *H' preserves all PPVs between nodes in S . Moreover, for any $0 < \rho < 1$ and any personalization vector \mathbf{r} such that $\mathbf{r}(j) = \rho$ for all $j \in V \setminus S$, $\mathbf{p}_r^G(i) = \mathbf{p}_r^{H'}(i)$, for all $i \in S$, where*

$$\mathbf{r}'(k) = \begin{cases} \mathbf{r}(k), & k \in S \\ \rho|V \setminus S|, & k = \mathbf{SOURCE} \\ 0, & k = \mathbf{SINK} \end{cases}$$

6. Experiments

Given a graph G and a subset S of its nodes we compare the subgraphs **WEIGHTED** and **UNWEIGHTED**

to the baseline INDUCED. As described in Section 5.2, WEIGHTED is obtained by retaining a subset of weighted edges from our PPV-preserving subgraph, and UNWEIGHTED removes the weights from the edges in WEIGHTED though edges to sink are still weighted. (Our rounding process ensures that WEIGHTED and UNWEIGHTED have the same number of edges as INDUCED.) We empirically determine how well each of these subgraphs preserves PPVs, and the usefulness of our subgraphs for applications like graph clustering.

DATA. We present results on three large real-world social networks. EPINIONS is a trust-based social network between the users of epinions.com website, with 75K nodes and 508K edges. SLASHDOT is a February 2009 snapshot of the friendship network between users of slashdot.org website, with 82K nodes and 948K edges. FACEBOOK is a snapshot, as of September 2005, of Facebook social network of the University of Oklahoma, with 17K nodes and 1.78M edges.

SAMPLING SUBSETS. Even though our algorithms can work on any subset of nodes, we show results for subsets generated according to prior work on graph sampling. We generated node samples using five techniques defined in (Leskovec & Faloutsos, 2006): (1) RN: select nodes uniformly at random; (2) RNN: select a node uniformly at random together with all its outneighbors; (3) RNE: select a node uniformly at random together with its random outneighbor; (4) RW: select a node uniformly at random and simulate a random walk with restart to the initial node (with probability 0.15), pick all traversed nodes; (5) RJ: identical to RW but with restart to a random node. For each generator, we experimented with subsets of size 100, 1000, and 10,000, though we only present results for subsets of size 1000 for lack of space. Results with other sizes are similar.

PRESERVING PPV. We compute the PPV from node i to node j for all $i, j \in S$, once using the entire graph and once restricted to just the subgraph under consideration. Figure 4 shows the Frobenius norm of the differences between these values, normalized by the Frobenius norm of the original PPV matrix. Observe that for all datasets and subset sampling methods, WEIGHTED improves over UNWEIGHTED by a factor of 3 on average, and both significantly outperform INDUCED. This agrees with our intuition: WEIGHTED is a closer approximation of our subgraphs which perfectly preserve PPVs; hence it performs better. UNWEIGHTED induces additional error since it loses the information encoded in the weights. The poor performance of INDUCED is in part due to the lack of a sink node. When most of a node’s neighbors in G are not

in S , the sink node allows us to better preserve PPV by absorbing the probability mass that was going to these neighbors. In fact, INDUCED and UNWEIGHTED share only 50% of their edges on average, showing that significant changes must be made to INDUCED to properly preserve PPVs.

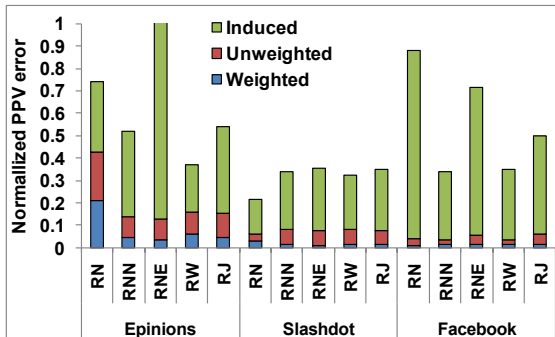


Figure 4. Preserving PPV: WEIGHTED is better than UNWEIGHTED by a factor of 3, and UNWEIGHTED in turn is better than INDUCED by a factor of 10 on average.

CLUSTERING ON SUBGRAPHS. We evaluate three hierarchical clustering algorithms on the INDUCED and UNWEIGHTED subgraphs: *Metis* (Karypis & Kumar, 1999), *Planted* (Condon & Karp, 1999), and *Local* (Andersen et al., 2008). The first two perform recursive graph bisection, while the third performs local clustering based on PageRank vectors. Note that the first two have no obvious relationship with random walks. For a sample S and a seed node $u \in S$, we order all the other nodes in S by their “clustering distance” from u , i.e., the length of the path in the hierarchy tree between the nodes. Intuitively, if the ordering in a subgraph is similar to that in G , then the clustering algorithm gets equivalent results when run on the subgraph or on G , and so any parameter tuning or model selection performed on the subgraph will also hold for G . To measure the similarity of two orderings we consider their prefixes of length $k = 10$ (that is the ten nodes closest to the seed) and compute the Kendall’s tau correlation coefficient. The experiment was run with each node in S as the seed node, and the results averaged. Fig. 5 shows, for each dataset and sampling method, the relative improvement of UNWEIGHTED over INDUCED (edges to sink were discarded, so UNWEIGHTED requires the same space as INDUCED here). We see that UNWEIGHTED yields a 50% improvement over INDUCED on average, demonstrating the usefulness of preserving PPVs, even approximately. More improvement is observed with *Metis* w.r.t. *Planted* because the latter uses *unnorm*-alized cuts; but also w.r.t. *Local* which uses PPVs. A possible explanation is that *Metis* uses a BFS and

(implicitly) shortest path distances, which are likely to be composed of a few important edges, whereas a PPV is an average over all paths between two nodes. Since rounding retains the most important edges, it is likely to preserve shortest paths (and hence, *Metis*) better than PPVs (and hence, *Local*). Finally, the RN subsets consistently yield the least improvement. This is because the induced subgraph for RN is typically extremely sparse (since nodes are picked randomly), and since UNWEIGHTED is allowed no more edges than INDUCED, there is little room for improvement.

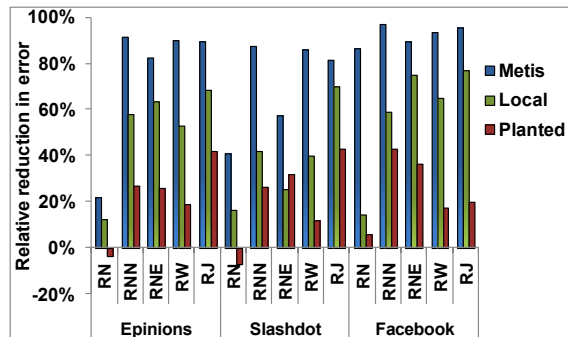


Figure 5. Reduction of error in preserving the clustering hierarchy, for UNWEIGHTED relative to INDUCED.

Conclusions

The PPVs of a graph are important for many applications such as clustering, classification, etc. Methods for representative subgraph construction typically consider only induced subgraphs and fail to preserve PPVs. We show how to build PPV-preserving subgraphs and how to sparsify them while still provably approximating the original PPVs. Experiments on three large real-world graphs show significant improvements of our method in maintaining PPVs, clustering accuracy, as well as in matching degree distributions and clustering coefficients (not shown for lack of space).

Acknowledgments. This work was supported in part by NSF Grant #0905645. We would also like to thank Purnamrita Sarkar for helpful discussions, and the reviewers for their insightful comments.

References

Andersen, R., Chung, F., and Lang, K. Local partitioning for directed graphs using pagerank. *Internet Math.*, 5(1-2):3–22, 2008.

Belkin, Mikhail and Niyogi, Partha. Semi-supervised learning on riemannian manifolds. *Machine Learning*, 56:209–239, 2004. ISSN 0885-6125.

Clauset, A., Shalizi, C., and Newman, M. Power-law

distributions in empirical data. *SIAM Review*’09.

- Condon, Anne and Karp, Richard M. Algorithms for graph partitioning on the planted partition model. *Random Structures & Algorithms*, 18:116–140, 1999.
- Garey, M. R., Johnson, D. S., and Stockmeyer, L. Some simplified np-complete problems. pp. 47–63, 1974. doi: <http://doi.acm.org/10.1145/800119.803884>.
- Gkantsidis, C., Mihail, M., and Saberi, A. Random walks in peer-to-peer networks: Algorithms and evaluation. *P2P Computing Systems*, 63(3), 2006.
- Gori, M. and Pucci, A. Itemrank: A random-walk based scoring algorithm for recommender engines. In *IJCAI*, pp. 2766–2771, 2007.
- Gyöngyi, Z., Garcia-Molina, H., and Pedersen, J. Combating web spam with trustrank. *VLDB*, 2004.
- Haveliwala, Taher H. Topic-sensitive pagerank: A context-sensitive ranking algorithm for web search. *IEEE Trans. Knowl. Data Eng.*, 15(4), 2003.
- Hübler, C., Kriegel, H.-P., Borgwardt, K., and Ghahramani, Z. Metropolis algorithms for representative subgraph sampling. In *ICDM*, 2008.
- Jeh, Glen and Widom, Jennifer. Scaling personalized web search. In *WWW*, pp. 271–279, 2003.
- Jia, Yuntao, Hoberock, Jared, Garland, Michael, and Hart, John. On the visualization of social and other scale-free networks. *IEEE Trans. on Visualization and Computer Graphics*, pp. 1285–1292, 2008.
- Kannan, Ravi, Vempala, Santosh, and Vetta, Adrian. On clusterings: Good, bad and spectral. *Journal of ACM*, 51, 2004.
- Karypis, George and Kumar, Vipin. A fast and highly quality multilevel scheme for partitioning irregular graphs. *Journal on Scientific Computing*, 1999.
- Leskovec, Jure and Faloutsos, Christos. Sampling from large graphs. In *KDD*, pp. 631–636, 2006.
- Maiya, Arun S. and Berger-Wolf, Tanya Y. Sampling community structure. In *WWW*, pp. 701–710, 2010.
- Meila, Marina and Shi, Jianbo. Learning segmentation by random walks. In *NIPS*, pp. 873–879, 2000.
- Moitra, Ankur. Approximation algorithms for multi-commodity type problems with guarantees independent of the graph size. In *FOCS*, pp. 3–12, 2009.
- Moser, Hannes and Thilikos, Dimitrios M. Parameterized complexity of finding regular induced subgraphs. *J. Discrete Algorithms*, 7(2):181–190, 2009.
- Ng, A., Zheng, A., and Jordan, M. Link analysis, eigenvectors, and stability. In *IJCAI*, 2001.
- Opsahl, T. and Panzarasa, P. Clustering in weighted

networks. *Social Networks*, 31(2):155–163, 2009.

Orponen, P., Schaeffer, S., and Gaytán, V. Locally computable approximations for spectral clustering and absorption times of random walks. *CoRR*, abs/0810.4061, 2008.

Sala, A., Cao, L., Wilson, C., Zablit, R., Zheng, H., and Zhao, B. Measurement-calibrated graph models for social network experiments. In *WWW*, 2010.

Sarlós, T., Benczúr, A., Csalogány, K., Fogaras, D., and Rácz, B. To randomize or not to randomize: space optimal summaries for hyperlink analysis. In *WWW*, pp. 297–306, 2006.

Szummer, M. and Jaakkola, T. Partially labeled classification with markov random walks. In *NIPS*, 2001.

van Dongen, Stijn Marinus. *Graph clustering by flow simulation*. PhD thesis, Univesity of Utrecht, 2000.

A. Induced subgraph

A.1. NP-hardness

For proving Thm. ??, we will use a reduction from the problem of finding a 2-regular induced subgraph of maximum size in *cubic*¹ graphs. This problem is known to be NP-hard for general graphs (Moser & Thilikos, 2009). The next two lemmas prove that the problem remains hard even when restricted to cubic graphs. We will start showing the hardness for graphs with nodes of degree at most three.

Lemma 11. *Finding a 2-regular induced subgraph of maximum size is NP-hard even in graphs of degree at most three.*

Proof. The reduction is from vertex cover which is known to be NP-hard even for cubic graphs (Garey et al., 1974). Given a cubic graph G of n nodes, we construct a graph G' as follows: we replace each node u in G with the gadget H_u (see fig. 6) and we distribute the 3 edges incident to u in G to the nodes u_1, u_2, u_3 in H_u so that u_1, u_2, u_3 have degree exactly 3 in the obtained graph G' (also note that every other node in G' has degree at most 3). We want to show that there is a vertex cover S of size k in G if and only if there is subset S' of nodes in G' with $|S'| = 6k + 3(n - k)$ such that removing S' from G' makes the graph 2-regular.

For one direction of the reduction, given a vertex cover S of size k we define S' in the following way: for each node $u \in S$, we add to S' the nodes $\{u_1, u_2, u_3, a_1, a_2, a_3\}$ of H_u , and for each node $u \notin S$, we add to S' the nodes $\{b_1, b_2, b_3\}$ of H_u . Note that

$|S'| = 6k + 3(n - k)$. Observe that if $u \in S$, then $H_u \setminus S'$ is 2-regular in $G' \setminus S'$ (it forms a cycle). If $u \notin S$ then all its neighbors must be in S : therefore, the “outside” edges incident to the nodes u_1, u_2, u_3 of H_u are not present in $G' \setminus S'$ which implies that $H_u \setminus S'$ is 2-regular (as a collection of three disjoint cycles).

For the other direction, let S' be any subset of nodes of G' of size $|S'| = 6k + 3(n - k)$ such that $G' \setminus S'$ is 2-regular. We make the following three claims for any H_u :

- (i) $S' \cap (\{c_1, c_2, c_3, d_1, d_2, d_3\} \cup P_{12} \cup P_{23} \cup P_{31}) = \emptyset$.
- (ii) If $u_i \in S'$ for some $i \in \{1, 2, 3\}$, then $S' \cap H_u = \{u_1, u_2, u_3, a_1, a_2, a_3\}$; in this case we say that H_u is *selected*.
- (iii) If u has neighbors v, w, z in G and H_u is not selected in G' , then $S' \cap H_u = \{b_1, b_2, b_3\}$ and all H_v, H_w, H_z are selected.

Note that the claim (iii) implies that $S = \{u \mid H_u \text{ is selected}\}$ is a vertex cover of G , and claims (ii)-(iii) implies that $|S| = k$. It remains to prove the claims.

To prove (i) assume $c_1 \in S'$. Then its neighbor in P_{12} has to be in S' too otherwise it would have degree one in $G' \setminus S'$, and recursively all the nodes in P_{12} have to be in S' . However, this is impossible since $|P_{12}| > 10 \cdot n$, while $|S'| = 6k - 3(n - k) < 10 \cdot n$. The same reasoning applies to the d 's nodes and to any of the nodes in the paths P_{12}, P_{23}, P_{31} .

For (ii), suppose wlog $u_1 \in S'$. We make the following chain of observations: $b_1 \notin S'$ (if not, d_1 has degree one in $G' \setminus S'$); $a_1 \in S'$ (if not, b_1 has degree 3 in $G' \setminus S'$); $a_2 \in S'$ (if not, c_1 has degree 3 in $G' \setminus S'$). Note that either b_2 or u_2 are in S' otherwise d_2 would have degree 3 in $G' \setminus S'$. If $b_2 \in S'$ then we know that all nodes $c_3, P_{31}, d_1, b_1, c_1, P_{12}, d_2, u_2$ are in $G' \setminus S'$ and are part of the same cycle which necessarily contains u_3 and b_3 , which is impossible since it would leave d_3 of degree 3. Therefore, we can conclude that $u_2 \in S'$. Applying the same argument, we obtain that a_3 and u_3 are the only other nodes in $S' \cap H_u$.

As for (iii), observe that claim (ii) implies that if $u_i \notin S'$ then $S' \cap \{u_1, u_2, u_3\} = \emptyset$. This in turn implies that b_1, b_2, b_3 are all in S' (if not, d_1, d_2, d_3 have degree 3). Also, note that no more nodes of H_u can be in S' , otherwise some node of H_u would have degree 1 in $G' \setminus S'$. Therefore, u_1, u_2, u_3 have degree 2 in $G' \setminus \{b_1, b_2, b_3\}$ only if all H_v, H_w, H_z are selected. \square

We now extend the hardness to cubic graphs.

¹A graph is cubic if all its nodes have degree exactly 3.

Lemma 12. *Finding a 2-regular induced subgraph of maximum size is NP-hard even in cubic graphs.*

Proof. We reduce from the problem of finding a 2-regular induced subgraph in graphs of degree at most 3 (Lemma 11). Let G be any such graph and construct a cubic graph G' from G as follows: to each node u append $(3 - \deg(u))$ copies of the gadget T depicted in fig. 7. Let t be the number of added gadgets. We claim that there is a 2-regular induced subgraph $G[U]$ with $|U| \geq k$ if and only if there is a 2-regular induced subgraph $G'[U']$ with $|U'| \geq k + 8t$: for one direction of the reduction, given U of size at least k and such that $G[U]$ is 2-regular, we obtain U' by adding to U the nodes in the two 4-cycles (a_1, b_1, c_1, d_1) and (a_2, b_2, c_2, d_2) of each gadget T added to G' . Clearly, $|U'| \geq k + 8t$ and $G'[U']$ is 2-regular. For the other direction, let U' be of size at least $k + 8t$ and such that $G'[U']$ is 2-regular. Then, every single cycle in $G'[U']$ is all contained either in G or in one of the added gadgets T . Let U be the nodes of U' contained in G . Then, since each gadget T can have at most 8 nodes inducing a 2-regular subgraph, we conclude that $|U| \geq |U'| - 8t \geq k$. \square

We are now ready to prove Thm ??.

Proof of Thm ??. We will prove the hardness of deciding whether there exists a *strict* subset $S \subsetneq V$ of size at least k such that $d(\tilde{\mathbf{p}}^G, \mathbf{p}^{G[S]}) = 0$. Then, the theorem easily follows. The proof is via a reduction from the problem of finding a 2-regular induced subgraph of maximum cardinality in cubic graphs which we proved to be NP-hard in Lemma 12. Given a cubic graph G on n nodes, we obtain another cubic graph G' by adding $t = 10 \cdot n$ disconnected tetrahedrons to G . We claim that there is a 2-regular induced subgraph $G[U]$ with $|U| \geq k$ if and only if there is a subset $S' \subsetneq V(G')$ of size at least $k + 3t$ such that $d(\tilde{\mathbf{p}}^{G'}, \mathbf{p}^{G'[S']}) = 0$.

For one direction of the reduction, given a set U of size at least k and such that $G[U]$ is 2-regular, we obtain S' by adding to U three nodes from each tetrahedron. Note that $G'[S']$ is 2-regular as well and $|S'| \geq k + 3t$ (and $S' \subsetneq V(G')$). Since $G'[S']$ and G' are both regular graphs, the PageRanks in each of the graphs are the same for all nodes — i.e., $\mathbf{p}^{G'}(i) = 1/|V(G')|$ for any $i \in V(G')$ and $\mathbf{p}^{G'[S']}(i) = 1/|V(G'[S'])| = 1/|S'|$ for any $i \in S'$, which implies $d(\tilde{\mathbf{p}}^{G'}, \mathbf{p}^{G'[S']}) = 0$. For the other direction of the reduction, consider any $S' \subsetneq V(G')$ of size at least $k + 3t$ such that $d(\tilde{\mathbf{p}}^{G'}, \mathbf{p}^{G'[S']}) = 0$. Since $d(\tilde{\mathbf{p}}^{G'}, \mathbf{p}^{G'[S']}) = 0$, then $\tilde{\mathbf{p}}^{G'} = \mathbf{p}^{G'[S']}$, or equivalently $\mathbf{p}^{G'[S']}(i) = 1/|S'|$ for every $i \in S'$. Given that all PageRanks in the graph $G'[S']$ are equal, it must be the case that the $G'[S']$ is regular (this fact follows from equation (1)).

Since $S' \subsetneq V(G')$ and G' is cubic, $G'[S']$ can be 0-regular (an independent set), 1-regular (a matching), or 2-regular. However, $G'[S']$ cannot be 0-regular because each tetrahedron can contribute only one node to S' , which means $|S'| < n + t < k + 3t$. Similarly, $G'[S']$ cannot be 1-regular because each tetrahedron can contribute only 2 nodes to S' , which means $|S'| < n + 2t < k + 3t$. We conclude that $G'[S']$ is 2-regular. Defining $U = S' \cap V(G)$, we have that $G[U]$ is 2-regular as well. Also, since each tetrahedron can contribute only 3 nodes to S' , it must be that $|U| \geq |S'| - 3t \geq k$ which concludes the proof. \square

A.2. Proof of Theorem 2

Proof. Rather than sampling a graph G from $G(n, p)$ and choosing a random subset S of G of size k , we consider the following equivalent 2-phase process: in the first phase, sample a graph $H = (S, E)$ from $G(k, p)$; for the second phase, let $H_0 = H$, and for $i > 0$ let H_i be the graph obtained by H_{i-1} by adding a new node u_i and adding each edge $\{u_i, v\}$ independently with probability p for all $v \in H_{i-1}$. By letting $G = H_{n-k}$, we can consider the graph H as the graph $G[S]$.

We claim that, for any $c > 0$, with probability at least $1 - O(k^{-c-1})$, every node in H has degree at most

$$g(c, p, k) := p \cdot (k - 1) + \sqrt{c \cdot p \cdot (k - 1) \cdot \log k}.$$

Observe that, for a fixed node $u \in S$, we have that $\deg_H(u) \sim \text{Bin}(k - 1, p)$. Therefore, from Chernoff's bound, we have that

$$\Pr[\deg_H(u) \geq (1 + \delta)p(k - 1)] < \exp(-p(k - 1)\delta^2/2).$$

Thus, with probability at least $1 - O(k^{-c})$, we have $\deg_H(u) \leq g(c, p, k)$. A union bound over the nodes in H concludes the proof of the claim. For the rest of the proof we condition on the event that every node in H has degree at most $g(c, p, k)$.

Let u^* be a node of maximum degree in H . To prove the theorem it is enough to show that, assuming $k = o(n/\log n)$, the degree of u^* in G is at most the average degree of G (i.e., $p \cdot (n - 1)$) with probability at least $1/2 - o(1)$. We will do so by finding a lower bound for the probability that the number of edges X added to u^* in the second phase is at most $p \cdot t - \epsilon \sqrt{p \cdot (1 - p) \cdot t}$, where $t = n - k$. Note that $X \sim \text{Bin}(t, p)$. By the central limit theorem we have that

$$\Pr \left[X \leq p \cdot t - \epsilon \sqrt{p \cdot (1 - p) \cdot t} \right] \geq \Phi(-\epsilon) - o(1),$$

where $\Phi(z)$ is the cumulative distribution of the normal distribution $N(0, 1)$. By Taylor expansion we have

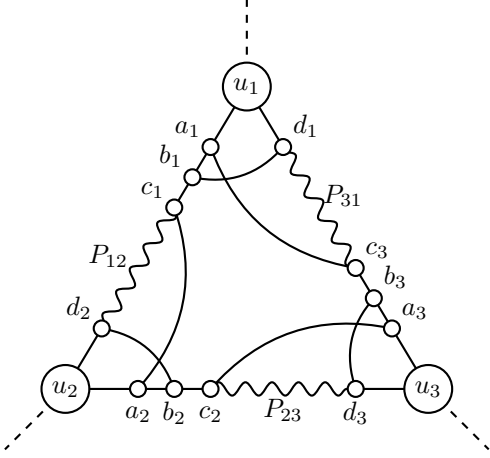


Figure 6. The gadget H_u : the paths P_{12} , P_{23} and P_{31} are of length more than $10 \cdot n$.

that $\Phi(\epsilon) \geq (1 - \epsilon(2/\sqrt{2\pi}))/2$. Therefore, with probability at least $1/2 - \epsilon/\sqrt{2\pi} - O(1/k^{c-1})$, the degree of u^* in G is

$$\deg_G(u^*) \leq g(c, p, k) + p \cdot t - \epsilon \sqrt{p \cdot (1-p) \cdot t}.$$

This implies that, with the same probability, $\deg_G(u^*) \leq p \cdot (k - 1 + t) = p \cdot (n - 1)$ when $t \geq \frac{c}{\epsilon^2(1-p)}(k-1) \log k$. The theorem follows since, for any $\epsilon > 0$, this condition on t is satisfied for $k = o(n/\log n)$ and n large enough. \square

B. Preserving PPVs

B.1. Constructing a PPVs-preserving subgraph

Proof of Theorem 4. The proof follows from Lemmas 5 and 6.

Proof of Lemma 5. The weights of the edges $e^*(x, y)$ are sums of non-negative components and are thus non-negative.

It remains to show that the weights on the edges to the sink ($e^*(x, \text{SINK})$) are non-negative. By construction of the algorithm and by the geometric series we have that $\sum_{y \neq z: w_G(z, y) > 0} w_H(e^*(x, y))$ is equal to

$$\begin{aligned} w_G(x, z) \frac{1 - \alpha}{1 - (1 - \alpha)w_G(z, z)} \sum_{y \neq z: w_G(z, y) > 0} w_G(z, y) \\ &= w_G(x, z) \frac{1 - \alpha}{1 - (1 - \alpha)w_G(z, z)} (1 - w_G(z, z)) \\ &= w_G(x, z) \frac{1 - (1 - \alpha)w_G(z, z) - \alpha}{1 - (1 - \alpha)w_G(z, z)} \leq w_G(x, z). \end{aligned}$$

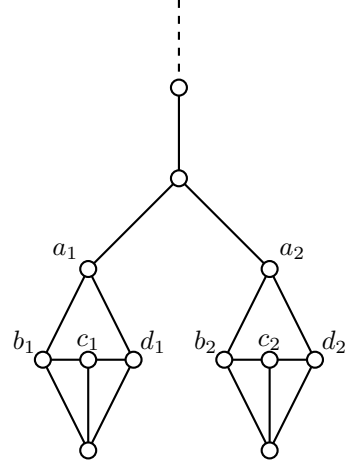


Figure 7. The gadget T .

The second equality holds since we assume the weights in G are normalized: $\sum_j w_G(i, j) = 1$ for all i .

The fact that weighted out-degree sums to 1 is easily seen algebraically. \square

Proof of Lemma 6. Let $z \in V \setminus S$ be the node removed by the algorithm so that the new graph is $H = ((V \setminus \{z\}) \cup \{\text{SINK}\}, E_H, w_H)$. Let x^- be an in-neighbor of z and x^+ be an out-neighbor of z such that $w_G(x^-, z) > 0$ and $w_G(z, x^+) > 0$ respectively. We denote by $e^*(x^-, x^+)$ the new directed edge (x^-, x^+) added by the algorithm in H .

Recall that we assume $\sum_{j \in V} w(i, j) = 1$ for every $i \in V$. By equations (2), the i -personalized PageRank $\mathbf{p}_i^G(j)$ with respect to node j in G is defined as follows:

$$\begin{aligned} \sum_{t=0}^{\infty} \alpha(1 - \alpha)^t \sum_{i=k_1, k_2, \dots, k_{t+1}=j} \prod_{l=1}^t w_G(k_l, k_{l+1}) \\ = \alpha \sum_{t=0}^{\infty} \sum_{i=k_1, k_2, \dots, k_{t+1}=j} \prod_{l=1}^t (1 - \alpha)w_G(k_l, k_{l+1}) \end{aligned}$$

Define $M_{i,j}^G$ as the set of all walks from i to j in G , and $h^G(s) = \prod_{r=1}^l (1 - \alpha)w_G(e_r)$, where s is a walk composed of the edges e_1, e_2, \dots, e_l . Then, we can write $\mathbf{p}_i^G(j) = \alpha \sum_{s \in M_{i,j}^G} h^G(s)$.

In order to prove the lemma, it suffices to show that, for all $i, j \in V \setminus \{z\}$, $\sum_{s \in M_{i,j}^G} h^G(s) = \sum_{s \in M_{i,j}^H} h^H(s)$. Henceforth, we interpret walks as sequences of edges. For two walks s_1 ending in i and s_2 starting from i , we define $s_1 \circ s_2$ to be the concatenation of the two. We say that a walk $s \in M_{i,j}^G$ is an (i, j) -basis (or simply basis

when clear from the context) if edge (z, z) does not appear in s , and we define $B_{i,j}^r$ to be the set of (i, j) -bases that pass through the node z exactly r times. For a basis $b \in B_{i,j}^r$ with $r \geq 1$, we let e_l^{b-} (resp. e_l^{b+}) be the l -th edge in b directed to (resp. directed out of) z and x_l^b, y_l^b be the corresponding nodes (i.e. $e_l^{b-} = (x_l^b, z)$ and $e_l^{b+} = (z, y_l^b)$). We define the closure of a basis $b \in B_{i,j}^r$ as $\text{cl}(b) = \bigcup_{t_1, \dots, t_r \geq 0} b^{(t_1, \dots, t_r)}$, where $b^{(t_1, \dots, t_r)}$ is the walk obtained from b by injecting t_l edges (z, z) between the edges e_l^{b-} and e_l^{b+} , for $1 \leq l \leq r$. For $b \in B_{i,j}^0$, we set $\text{cl}(b) = \{b\}$. Observe that for any walk $s \in M_{i,j}^G$, there exists exactly one basis $b \in \bigcup_{r \geq 0} B_{i,j}^r$ such that $s \in \text{cl}(b)$.

We define a mapping ϕ from bases in G to walks in H as follows: for $b \in B_{i,j}^r$, $\phi(b)$ is obtained from b by replacing the edges $e_l^{b-} \circ e_l^{b+}$ in b with the new edge $e^*(x_l^b, y_l^b)$ added in the graph H , for $1 \leq l \leq r$ (note that $\phi(b) = b$ for $b \in B_{i,j}^0$). Observe that for every walk $s \in M_{i,j}^H$, there exists exactly one basis $b \in B_{i,j}^r$ (for some r) such that $s = \phi(b)$, so the function ϕ is bijective.

We claim that for every $i, j \in S$, for every $r \geq 0$ and $b \in B_{i,j}^r$, $\sum_{s \in \text{cl}(b)} h^G(s) = h^H(\phi(b))$. Before proving the claim we show how the lemma follows from it:

$$\begin{aligned} \sum_{s \in M_{i,j}^G} h^G(s) &= \sum_{r=0}^{\infty} \sum_{b \in B_{i,j}^r} \sum_{s \in \text{cl}(b)} h^G(s) \\ &= \sum_{r=0}^{\infty} \sum_{b \in B_{i,j}^r} h^H(\phi(b)) \end{aligned}$$

which is equivalent to $\sum_{s \in M_{i,j}^H} h^H(s)$ since $\phi_{i,j}$ is bijective.

It remains to prove the claim. We use induction on r . For $r = 0$, the equality holds since for every $b \in B_{i,j}^0$, we have $\text{cl}(b) = \{b\}$, $\phi(b) = b$, and $h^G(b) = h^H(b)$ because b uses only edges existing in G . Now consider any $r > 1$ and suppose the claim holds for $r-1$. Consider any $b \in B_{i,j}^r$ and let b' (resp. b'') be the sequence of edges of b before e_1^{b-} (resp. after e_1^{b+}). We have that $b'' \in B_{i,j}^{r-1}$. Note that $\text{cl}(b)$ contains exactly the walks of the form $b' \circ e_1^{b-} \circ (z, z) \circ \dots \circ (z, z) \circ e_1^{b+} \circ s$, where $s \in \text{cl}(b'')$ and the sequence $(z, z) \circ \dots \circ (z, z)$ is of length zero or more. We observe that $b' = \phi(b')$ and $h^G(b') = h^H(b')$. Also, by construction of the algorithm for any l we have

$$h^H(e^*(x_l^b, y_l^b)) = h^G(e_l^{b-}) \cdot h^G(e_l^{b+}) \sum_{t \geq 0} [h^G(z, z)]^t.$$

Therefore, $\sum_{s \in \text{cl}(b)} h^G(s)$ is equal to

$$\begin{aligned} &h^G(b') \cdot h^G(e_1^{b-}) \sum_{t \geq 0} [h^G(z, z)]^t \cdot h^G(e_1^{b+}) \sum_{s \in \text{cl}(b'')} h^G(s) \\ &= h^G(b') \cdot h^H(e^*(x_1^b, y_1^b)) \sum_{s \in \text{cl}(b'')} h^G(s) \\ &= h^H(\phi(b')) \cdot h^H(\phi(e_1^{b-} \circ e_1^{b+})) \sum_{s \in \text{cl}(b'')} h^G(s) \\ &= h^H(\phi(b' \circ e_1^{b-} \circ e_1^{b+})) \sum_{s \in \text{cl}(b'')} h^G(s) \\ &= h^H(\phi(b' \circ e_1^{b-} \circ e_1^{b+})) \cdot h^H(\phi(b'')) = h^H(\phi(b)), \end{aligned}$$

where the last two steps follow by induction and by the definitions of h^H and ϕ respectively. \square

Proof of Lemma 7. From Theorem 4, we know that there exists some setting of weights W_H such that (a) $\sum_j W_H(i, j) = 1$ for all i , and (b) the PPV matrix is identical to the desired matrix P^H . Let $P^H(j)$ be the j -th column of P^H . From equation (1), we must have:

$$\begin{aligned} P^H(j) &= \alpha e_j + (1 - \alpha) i W_H^t P^H(j) \\ \Rightarrow [(1 - \alpha) W_H^t - I] P^H(j) &= -\alpha e_j \\ \Rightarrow [(1 - \alpha) W_H^t - I] P^H &= -\alpha I \quad (3) \\ \Rightarrow W_H^t &= \frac{1}{1 - \alpha} (I - \alpha (P^H)^{-1}) \end{aligned}$$

where the inverse in the last step exists because each factor in the LHS of Eq. 3 must be full rank (since the RHS is full rank). \square

B.2. Rounding the Weighted Matrix

Proof of Theorem 8. The proof is similar to that of Theorem 3 of (Ng et al., 2001). Given a node $j \in S$, we construct a coupled Markov chain $\{(X_t, Y_t) | t \geq 0\}$ on the set of nodes $S \cup \{\text{SINK}\}$ as follows. Initially we set $X_0 = Y_0 = j$. At any step t , with probability α , we “reset” both chains to node j , otherwise we draw a node ℓ according to the weight matrix w_H , and set $X_{t+1} = \ell$. Now, if $X_t = Y_t = k$ and the edge (k, ℓ) was not rounded by WEIGHTED, then we also set $Y_{t+1} = \ell$. If either condition is not satisfied, we draw a node ℓ' according to the rounded weight matrix w_R and set $Y_{t+1} = \ell'$.

Now, we compute the probability $d_{t+1} = \Pr(X_{t+1} \neq Y_{t+1})$. This event occurs only if there was no reset at

step $t + 1$. Thus,

$$\begin{aligned}
 \frac{d_{t+1}}{1-\alpha} &= \Pr(X_{t+1} \neq Y_{t+1} | \text{no reset at } t+1) \\
 &\leq \Pr(X_t \neq Y_t) \\
 &\quad + \sum_k \Pr(X_t = Y_t = k) \Pr((k, X_{t+1}) \text{ was rounded}) \\
 &\leq \Pr(X_t \neq Y_t) \\
 &\quad + \sum_k \Pr(X_t = Y_t = k) \sum_\ell \Pr(X_{t+1} = \ell, w_H(k, \ell) < \tau) \\
 &\leq d_t + \sum_k \Pr(X_t = Y_t = k) \sum_{\{\ell | w_H(k, \ell) < \tau\}} \Pr(X_{t+1} = \ell) \\
 &\leq d_t + \sum_k \Pr(X_t = Y_t = k) \sum_{\{\ell | w_H(k, \ell) < \tau\}} w_H(k, \ell) \\
 &\leq d_t + \sum_k \Pr(X_t = Y_t = k) \tau(|S| + 1) \\
 &\leq d_t + \tau(|S| + 1)
 \end{aligned}$$

Using this recurrence and the fact that $d_0 = 0$, we find that $d_\infty \leq \tau|S|(1-\alpha)/\alpha$. Now, the stationary distributions of the X and Y Markov chains are just the PPV vectors $P^H(j)$ and $P^R(j)$. By the Coupling Lemma, $\|P^H(j) - P^R(j)\|_1 \leq 2d_\infty$, which yields the theorem statement.

B.3. Preserving Global Pageranks

Proof of Corollary 9. Using the linearity of equation (1), we decompose $\mathbf{p}_r^G(i)$ for any $i \in S$ into two components: one summarizes all contributions from nodes in S and another one from all nodes in $V \setminus S$:

$$\mathbf{p}_r^G(i) = \sum_{j \in V} \mathbf{r}(j) \mathbf{p}_j^G(i) = \sum_{j \in S} \mathbf{r}(j) \mathbf{p}_j^G(i) + \sum_{j \in V \setminus S} \mathbf{r}(j) \mathbf{p}_j^G(i).$$

The first sum can be simply computed as $\sum_{j \in S} \mathbf{r}(j) \mathbf{p}_j^G(i) = \sum_{j \in S} \mathbf{r}(j) \mathbf{p}_j^H(i)$ by Lemma 6. As for the second sum, by the assumption on \mathbf{r} we have that

$$\sum_{j \in V \setminus S} \mathbf{r}(j) \mathbf{p}_j^G(i) = \mathbf{r}(j') \sum_{j \in V \setminus S} \mathbf{p}_j^G(i) = \mathbf{r}(j') \mathbf{v}(i).$$

for any $j' \in V \setminus S$. The corollary follows since $\mathbf{r}(j')$ is a constant: $\mathbf{r}(j') = (1 - \sum_{j \in S} \mathbf{r}(j))/|V \setminus S|$. \square

Proof of Thm. 10. For every $i, j \in S$, no walk from i to j passes through SOURCE in H' , therefore PPVs between nodes in S are still preserved.

For the remaining part of the theorem, we recall that the personalized PageRank corresponding to any vector \mathbf{r} can be computed as a linear combination

of the personalized PageRank vectors of individual nodes (Haveliwala, 2003; Jeh & Widom, 2003).

$$\mathbf{p}_r = \sum_{j \in V} \mathbf{r}(j) \cdot \mathbf{p}_j. \quad (4)$$

By equation (4), we have that

$$\begin{aligned}
 \mathbf{p}_r^G(i) &= \sum_{j \in S} \mathbf{r}(j) \mathbf{p}_j^G(i) + \sum_{j \in V \setminus S} \mathbf{r}(j) \mathbf{p}_j^G(i), \\
 \mathbf{p}_r^{H'}(i) &= \sum_{j \in S} \mathbf{r}'(j) \mathbf{p}_j^{H'}(i) + \sum_{j \in \{\text{SOURCE}, \text{SINK}\}} \mathbf{r}'(j) \mathbf{p}_j^{H'}(i).
 \end{aligned}$$

By the definition of \mathbf{r}' and the fact that $p_j^G(i) = p_j^{H'}(i)$ for $i, j \in S$, it suffices to show that

$$|V \setminus S| \cdot \mathbf{p}_{\text{SOURCE}}^{H'}(i) = \sum_{j \in V \setminus S} \mathbf{p}_j^G(i).$$

As in the proof of Lemma 6, we define $M_{j,i}^G$ to be the set of walks from j to i in G , and for $s \in M_{j,i}^G$, we define $h^G(s) = \prod_{e \in s} (1-\alpha)w_G(e)$. Therefore, by equations (2) and (3), $\mathbf{p}_j^G(i) = \alpha \sum_{s \in M_{j,i}^G} h^G(s)$. Since $j \in V \setminus S$ and $i \in S$, each walk in $M_{j,i}^G$ has an edge (k, l) where $k \in V \setminus S$ and $l \in S$, and first part of the walk does not go through any node in S before reaching (k, l) . Thus, we can rewrite $\mathbf{p}_j^G(i)$ as

$$\mathbf{p}_j^G(i) = \alpha \sum_{k \in V \setminus S} \sum_{l \in S} \sum_{\substack{s_0 \in M_{j,k}^G \\ s_0 \cap S = \emptyset}} \sum_{s_1 \in M_{l,i}^G} h^G(s_0) \cdot h^G(k, l) \cdot h^G(s_1).$$

Observe that $\alpha \sum_{s_0 \in M_{j,k}^G, s_0 \cap S = \emptyset} h^G(s_0) = \mathbf{p}_l^{G'}(i)$, and that $\alpha \sum_{s_1 \in M_{l,i}^G} h^G(s_1) = \mathbf{p}_l^G(i)$. Therefore, by rearranging the terms, we get

$$\mathbf{p}_j^G(i) = \frac{1}{\alpha} \sum_{k \in V \setminus S} \mathbf{p}_j^{G'}(k) \sum_{l \in S} \mathbf{p}_l^G(i) \cdot h^G(k, l).$$

By summing over $j \in V \setminus S$ and rearranging the sums we obtain

$$\sum_{j \in V \setminus S} \mathbf{p}_j^G(i) = (1-\alpha) \sum_{l \in S} \mathbf{p}_l^G(i) \cdot \frac{1}{\alpha} \sum_{j, k \in V \setminus S} \mathbf{p}_j^{G'}(k) \cdot w_G(k, l).$$

Now, since $\mathbf{p}_{\text{SINK}}^{H'}(i) = 0$ for all $i \in S$, and $\mathbf{p}_l^{H'} = \mathbf{p}_l^G$ for all $l \in S$, the decomposition theorem (Jeh & Widom, 2003) implies

$$\mathbf{p}_{\text{SOURCE}}^{H'}(i) = (1-\alpha) \sum_{l \in S} \mathbf{p}_l^G(i) \cdot w_{H'}(\text{SOURCE}, l)$$

The theorem follows by substituting $w_{H'}(\text{SOURCE}, l)$. \square

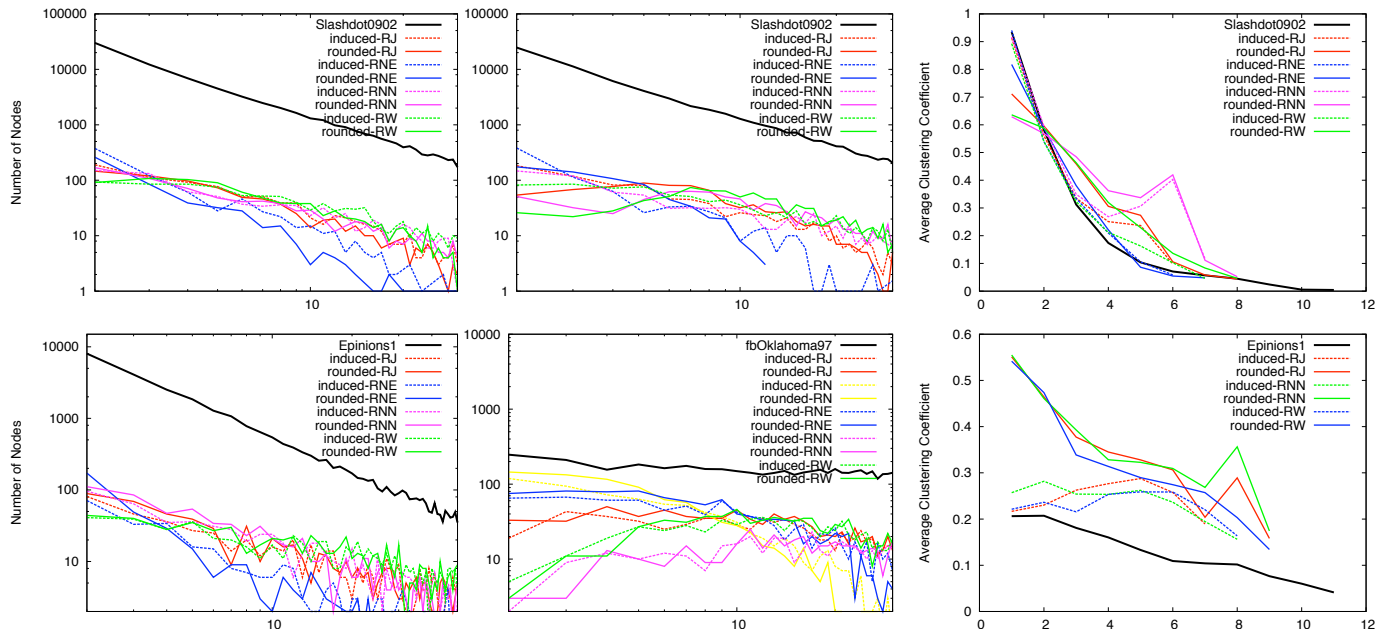


Figure 8. In-degree (1-st column), out-degree (2-nd column) distributions, and binned clustering coefficients (3-rd column) of original graphs and subgraphs.

C. Experiments

GRAPH STRUCTURE. As another application, we show that our method can be used to augment the induced graph with “long-distance” edges in order to preserve high-level graph structure. For ease of visualization, we used an HB/bcsstk26 graph², and made it unweighted and undirected by discarding edge weights and directions, obtaining 1,922 nodes and 14K undirected edges. The graph, visualized using Graphviz³, is shown in Figure 1. We sampled 200 nodes ($\approx 10\%$) using the RJ method. The sampled nodes are shown by black points. The induced subgraph on the sample nodes is shown in Figure 2. Obviously, the induced subgraph preserves only connectivity between the immediate neighbors, failing to preserve the high-level graph structure and leading to many disconnected components. Clearly, this induced subgraph is useless for any further processing relying on graph manifold structure, such as classification with random walks (Szummer & Jaakkola, 2001), or learning on manifolds (Belkin & Niyogi, 2004). The rounded graph (the sink node and all edge weights discarded, average node degree is as in the original graph) is shown in Figure 3. Observe that all but two components are merged into a single connected component. The “bridging” edges are created between components that were most tightly connected (according to PPV) in the original graph. We conclude that our rounding algorithm can

be used for preserving graph structure for small sample sizes where induced subgraphs fail.

DEGREE DISTRIBUTIONS AND CLUSTERING COEFFICIENTS. In this section, we consider the fraction of nodes with a given out- or in-degree in the subgraph. Degree distributions in real-world graphs typically follow power laws (Clauset et al.). The first two plots in Figure 8 (from left to right) depict the in- and out-degree distributions of all generated subgraphs for SLASHDOT. The degree distributions of the original SLASHDOT graph are depicted in these plots as well. We observe that for any sampling method, the distributions of UNWEIGHTED and the corresponding INDUCED are very similar to each other, and have about the same shape (slope) as the true degree distribution (in general, RJ and RW seem to approximate the true shape slightly better)⁴. Similar behaviors can be observed for the degree distributions of the other datasets: for presentation purposes we only show the in-degree distributions for EPINIONS and the (undirected-) degree distributions for FACEBOOK (first two plots of second row of Figure 8). We conclude that UNWEIGHTED preserves degree distributions as much as INDUCED, while outperforming INDUCED in terms of PPVs.

The clustering coefficient of a directed graph measures the probability that for any pair of directed edges (u, v)

⁴We ignore weights, so WEIGHTED and UNWEIGHTED have the same degree distributions and clustering coefficients.

²www.cise.ufl.edu/research/sparse/matrices

³<http://www.graphviz.org/>

and (v, w) sharing a node v , there is a third edge (u, w) (Opsahl & Panzarasa, 2009). This is a measure of transitivity of link structure in the graph, and is typically high for social networks. Instead of presenting just one number for the entire graph, we plot the clustering coefficient as a function of degree, after grouping degrees into bins of logarithmic sizes so that each bin has enough nodes to compute the clustering coefficient for that bin reliably. We only measure the clustering coefficient on unweighted graphs. The last column of Figure 8 depicts the binned distribution of clustering coefficients of INDUCED, UNWEIGHTED, as well as the true clustering coefficients, for SLASHDOT and EPINIONS. We observe that, for SLASHDOT, all subgraphs (except for the RNN sampling method) closely approximate the true clustering coefficient. In addition, for any given sampling method, UNWEIGHTED and INDUCED are very similar to each other. For EPINIONS, UNWEIGHTED is still substantially preserving the true shape (except for the first bins), while INDUCED follows a different trend by increasing and then decreasing. One possible explanation of the high clustering coefficient of nodes of low degree in UNWEIGHTED is the fact that these nodes might have outgoing edges to their “important” nodes in UNWEIGHTED that compensate for the paths that were present in the original graph. In INDUCED this rectification is not present and only a few edges are induced between low degree nodes, therefore the clustering coefficient is low.