

Semi-Supervised Classification through the Bag-of-Paths Group Betweenness

Bertrand Lebichot, Ilkka Kivimäki, Kevin François & Marco Saerens, *Member, IEEE*

Abstract—This paper introduces a novel, well-founded, betweenness measure, called the Bag-of-Paths (BoP) betweenness, as well as its extension, the BoP group betweenness, to tackle semi-supervised classification problems on weighted directed graphs. The objective of semi-supervised classification is to assign a label to unlabeled nodes using the whole topology of the graph and the labeled nodes at our disposal. The BoP betweenness relies on a bag-of-paths framework [1] assigning a Boltzmann distribution on the set of all possible paths through the network such that long (high-cost) paths have a low probability of being picked from the bag, while short (low-cost) paths have a high probability of being picked. Within that context, the BoP betweenness of node j is defined as the sum of the a posteriori probabilities that node j lies in-between two arbitrary nodes i, k , when picking a path starting in i and ending in k . Intuitively, a node typically receives a high betweenness if it has a large probability of appearing on paths connecting two arbitrary nodes of the network. This quantity can be computed in closed form by inverting a $n \times n$ matrix where n is the number of nodes. For the group betweenness, the paths are constrained to start and end in nodes within the same class, therefore defining a within-class group betweenness for each class. Unlabeled nodes are then classified according to the class showing the highest group betweenness. Experiments on various real-world data sets show that BoP group betweenness performs competitively compared to all the tested state-of-the-art methods [2]–[5]. The benefit of the BoP betweenness is particularly noticeable when only a few labeled nodes are available.

Index Terms—Graph and network analysis, network data, graph mining, betweenness centrality, kernels on graphs, semi-supervised classification.

I. INTRODUCTION

THE goal of a classification task is to automatically assign data to predefined classes. Traditional pattern recognition, machine learning or data mining classification methods require large amounts of labeled training instances, which are often difficult to obtain. The effort required to label the data can be reduced using, for example, semi-supervised learning methods. This name comes from the fact that the used data is a mixture of data used for supervised and unsupervised learning (see, e.g., [6] for an introduction). In short, semi-supervised learning methods learn from both labeled and unlabeled instances. This allows to reduce the amount of labeled instances needed to achieve the same level of classification accuracy.

Graph-based semi-supervised classification has received a growing focus in recent years. The problem can be described

as follows: given an input graph with some nodes labeled, the goal is to predict the missing node labels. This problem has numerous applications such as classification of individuals in social networks, categorization of linked documents (e.g. patents or scientific papers), or protein function prediction, to name a few. In this kind of application (as in many others), unlabeled data are usually available in large quantities and are easy to collect: friendship links can be recorded on Facebook, text documents can be crawled from the internet and DNA sequences of proteins are readily available from gene databases. Given a relatively small labeled data set and a large unlabeled data set, semi-supervised algorithms can infer useful information from both sources.

This paper tackles this problem of semi-supervised classification within the **bag-of-paths (BoP) framework**. This framework was originally introduced in the context of distance computation on graphs [1], capturing its global structure with, as building block, network paths. This same framework was previously used in [7] for defining a covariance kernel on a graph. More precisely, we assume a weighted directed graph or network G where a transition cost is associated to each arc. We further consider a bag containing all the possible paths (or walks) between pairs of nodes in G . Then, a Boltzmann distribution, depending on a temperature parameter T , is defined on the set of paths such that long (high-cost) paths have a low probability of being picked from the bag, while short (low-cost) paths have a high probability of being picked. In this framework, the **BoP probabilities**, $P(s = i, e = j)$, of sampling a path starting in node i and ending in node j can easily be computed in closed form by a simple $n \times n$ matrix inversion, where n is the number of nodes.

Within this context, a betweenness measure quantifying to which extent a node j is in-between two nodes i and k is defined. More precisely, the **BoP betweenness**, $\text{bet}_j = \sum_{i=1}^n \sum_{k=1}^n P(\text{int} = j | s = i, e = k)$, of a node j of interest is defined quite naturally as the sum of the a posteriori probabilities that node j (intermediate node) lies on a path between the two nodes i and k sampled from the graph bag-of-paths Boltzmann distribution. Intuitively, a node receives a high betweenness if it has a large probability of appearing on paths connecting two arbitrary nodes of the network.

For the **group betweenness**, the paths are constrained to start and end in nodes of the same class, therefore defining a group betweenness between classes, $\text{gbet}_j(C_i, C_k) = P(\text{int} = j | s \in C_i, e \in C_k)$. Unlabeled nodes are then classified according to the class showing the highest group betweenness when starting and ending within the same class.

In summary, this work has three main contributions:

The authors are with Universite Catholique de Louvain, ICTEAM & LSM (e-mail: bertrand.lebichot@uclouvain.be).

This work was partially supported by the Elis-IT project funded by the “Région wallonne”. We thank this institution for giving us the opportunity to conduct both fundamental and applied research.

- It develops both a betweenness measure and a group betweenness measure from a well-founded theoretical framework, the bag-of-paths framework. These two measures can be easily computed in closed form.
- This group betweenness measure provides a new algorithm for graph-based semi-supervised classification.
- It assesses the accuracy of the proposed algorithm on thirteen standard data sets and compares it to state-of-the-art techniques. The obtained performances are competitive with the other graph-based semi-supervised techniques.

The main drawback of the proposed method is that it requires a matrix inversion and therefore does not scale to very large graphs. An approximate method relying on bounded walks (as in [8]) will be investigated and is left for further work.

In this paper, the **BoP classifier** (or just BoP) will refer to the semi-supervised classification algorithm based on the bag-of-paths group betweenness developed in Section V.

The paper is organized as follows. Section II introduces background and notations, mainly the bag-of-paths and the bag-of-hitting-paths models. Then, related work in semi-supervised classification is discussed in Section III. The bag-of-paths betweenness and group betweenness centralities are introduced in Section IV. This enables us to derive the BoP classifier in Section V. Then experiments involving the BoP classifier and classifiers discussed in the related work section will be performed in Section VI. Results and discussions of those experiments can be found in Section VI-C. Finally, Section VII concludes this paper and opens a reflection for further work.

II. BACKGROUND AND NOTATION

This section aims to introduce the theoretical background and notation used in this paper. Furthermore, in order to understand the bag-of-paths betweenness of Section IV, the bag-of-paths and bag-of-hitting-paths frameworks [1] need to be reviewed first. Thus, graph-based semi-supervised classification will be discussed in Subsection II-A, then the bag-of-paths model will be introduced in Subsection II-B and, finally, the bag-of-hitting-paths model will be briefly described in Subsection II-C.

A. Graph-based semi-supervised classification

Consider a weighted directed graph or network, G , strongly connected with a set of n nodes \mathcal{V} (or vertices) and a set of edges \mathcal{E} (or arcs, links). The **adjacency matrix** of the graph, containing non-negative affinities between nodes, is denoted as \mathbf{A} , with elements $a_{ij} \geq 0$.

Moreover, to each edge between node i and j is associated a non-negative number $c_{ij} \geq 0$. This number represents the **immediate cost of transition** from node i to j . If there is no link between i and j , the cost is assumed to take a large value, denoted by $c_{ij} = \infty$. The **cost matrix** \mathbf{C} is an $n \times n$ matrix containing the c_{ij} as elements. Costs are set independently of the adjacency matrix – they are quantifying the cost of a transition according to the problem at hand. Costs can, e.g., be set in function of some properties, or features, of the nodes or the arcs in order to bias the probability distribution of choosing

a path. In the case of a social network, we may, for instance, want to bias the paths in function of the education level of the persons, therefore favoring paths visiting highly educated persons (see [1] for details). Now, if there is no reason to introduce a cost, we simply set $c_{ij} = 1$ (paths are penalized by their length) or $c_{ij} = 1/a_{ij}$ (in this case, a_{ij} is viewed as a conductance and c_{ij} as a resistance) – this last setting will be used in the experimental section.

When tackling the semi-supervised classification problem, we will also consider a set of classes, $\{C_k\}_{k=1}^m$, with the number of classes equal to m . Each node is assumed to belong to at most one class since the class label can also be unknown. To represent the class memberships, a $n \times m$ -dimensional indicator matrix, \mathbf{Y} , is used. On each of its rows, it contains, as entries, a 1 in column c when the corresponding node belongs to class c , and 0 otherwise (m zeros on line i if the node i is unlabeled). The c -th column of \mathbf{Y} will be denoted \mathbf{y}^c and contains the binary memberships of the nodes to class c .

Moreover, a **natural random walk** on G is defined in the standard way. In node i , the random walker chooses the next edge to follow according to reference transition probabilities

$$p_{ij}^{\text{ref}} = \frac{a_{ij}}{\sum_{j'=1}^n a_{ij'}} \quad (1)$$

representing the probability of jumping from node i to node $j \in \text{Succ}(i)$, the set of successor nodes of i . The corresponding transition probabilities matrix will be denoted as \mathbf{P}^{ref} . In other words, the random walker chooses to follow an edge with a probability proportional to the affinity (apart from the sum-to-one normalization), therefore favoring edges associated to a large affinity. The matrix \mathbf{P}^{ref} , containing the p_{ij}^{ref} , is stochastic and is called the **reference transition matrix**.

B. The bag-of-paths framework

This framework was recently introduced in [1] (see also [7]) for computing distances on graphs. In order to make the paper as self-contained as possible, we will briefly review the whole framework in this section (see [1] for details) and then use it in order to define a new betweenness measure, which is the main contribution of the paper. The bag-of-paths (BoP) model can be considered as a motif-based model [9], [10] using, as building block, paths of the network. In the next subsection, hitting paths will be used instead, as motifs.

A path φ (sometimes called a walk) is a sequence of transitions to adjacent nodes on G (loops are allowed), initiated from a starting node s , and stopping in an ending node e . If we want to emphasise on those starting and ending nodes, we will use φ_{se} instead of φ .

The BoP framework is based on the probability of picking a path $i \rightsquigarrow j$ starting at a node i and ending in a node j from a virtual bag containing all possible paths in the network [1]. Let us define \mathcal{P}_{ij} as the set of all possible paths connecting node i to node j , including loops. We further define the set of all paths through the graph as $\mathcal{P} = \bigcup_{i,j=1}^n \mathcal{P}_{ij}$. The **total cost** of a path φ , $\tilde{c}(\varphi)$, is defined as the sum of the individual transition costs c_{ij} along φ .

The potentially infinite set of paths in the graph is enumerated and a probability distribution is assigned to the set of individual paths \mathcal{P} . This probability distribution on the set \mathcal{P} , represents the probability of picking a path $\varphi \in \mathcal{P}$ in the bag, and is defined as the probability distribution \mathbf{P} minimizing the total expected cost, $\mathbb{E}[\tilde{c}(\varphi)]$, among all the distributions having a fixed relative entropy J_0 with respect to a reference distribution, for instance a natural random walk on the graph [1]. This choice naturally defines a probability distribution on the set of paths such that long (high cost) paths occur with a low probability while short (low cost) paths occur with a high probability. In other words, we are seeking path probabilities, $\mathbf{P}(\varphi)$, $\varphi \in \mathcal{P}$, minimizing the total expected cost subject to a constant relative entropy constraint:

$$\left\{ \begin{array}{l} \underset{\mathbf{P}(\varphi)}{\text{minimize}} \quad \sum_{\varphi \in \mathcal{P}} \mathbf{P}(\varphi) \tilde{c}(\varphi) \\ \text{subject to} \quad \sum_{\varphi \in \mathcal{P}} \mathbf{P}(\varphi) \ln(\mathbf{P}(\varphi)/\tilde{\pi}^{\text{ref}}(\varphi)) = J_0 \\ \sum_{\varphi \in \mathcal{P}} \mathbf{P}(\varphi) = 1 \end{array} \right. \quad (2)$$

where $\tilde{\pi}^{\text{ref}}$ represents the probability of following the path φ when walking according to the natural random walk reference distribution. Thus $\tilde{\pi}^{\text{ref}}$ is the product of the transition probabilities p_{ij}^{ref} along the path φ . Here, $J_0 > 0$ is provided a priori by the user, according to the desired degree of randomness, or exploration, he is willing to concede.

The result of the minimization [1] is a **Boltzmann probability distribution**:

$$\mathbf{P}(\varphi) = \frac{\tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)]}{\sum_{\varphi' \in \mathcal{P}} \tilde{\pi}^{\text{ref}}(\varphi') \exp[-\theta \tilde{c}(\varphi')]} \quad (3)$$

where $\theta = 1/T$ plays the role of an inverse temperature. As expected, short paths φ (having a low $\tilde{c}(\varphi)$) are favored in that they have a larger probability of being chosen. Moreover, from Equation (3), we clearly observe that when $\theta \rightarrow 0^+$, the paths probabilities reduce to the probabilities generated by the natural random walk on the graph. In this case, $J_0 \rightarrow 0$ and paths are chosen according to their likelihood in a natural random walk. On the other hand, when θ is large, the probability distribution defined by Equation (3) is biased towards short paths (paths shortest ones are more likely). Notice that, in the sequel, it will be assumed that the user provides the value of the parameter θ instead of J_0 , with $\theta > 0$.

The **bag-of-paths probability** [1] is now defined as the quantity $\mathbf{P}(s = i, e = j)$ on the set of (starting, ending) nodes of the paths. It corresponds to the probability of drawing a path starting in node i and ending in node j from the virtual bag-of-paths:

$$\mathbf{P}(s = i, e = j) = \frac{\sum_{\varphi \in \mathcal{P}_{ij}} \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)]}{\sum_{\varphi' \in \mathcal{P}} \tilde{\pi}^{\text{ref}}(\varphi') \exp[-\theta \tilde{c}(\varphi')]} \quad (4)$$

where \mathcal{P}_{ij} is the set of paths connecting the starting node i to the ending node j .

Let us derive the analytical closed form of this expression. To this end, we start from the cost matrix, \mathbf{C} , from which we build a new matrix, \mathbf{W} , as

$$\mathbf{W} = \mathbf{P}^{\text{ref}} \circ \exp[-\theta \mathbf{C}], \quad (5)$$

where \mathbf{P}^{ref} is the reference transition matrix containing the p_{ij}^{ref} , the exponential function is taken elementwise and \circ is the elementwise multiplication (Hadamard product). The entries of \mathbf{W} are therefore $w_{ij} = p_{ij}^{\text{ref}} \exp[-\theta c_{ij}]$.

It is shown in [1] that the numerator of Equation (4) is

$$\sum_{\varphi \in \mathcal{P}_{ij}} \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)] = \sum_{t=0}^{\infty} [\mathbf{W}^t]_{ij} = \left[\sum_{t=0}^{\infty} \mathbf{W}^t \right]_{ij} \quad (6)$$

where, by convention, zero-length paths are taken into account and are associated to a zero cost. Computing the series of powers of \mathbf{W} provides

$$\sum_{t=0}^{\infty} \mathbf{W}^t = (\mathbf{I} - \mathbf{W})^{-1} = \mathbf{Z}, \quad (7)$$

which converges if the spectral radius of \mathbf{W} is less than 1, $\rho(\mathbf{W}) < 1$. Since the matrix \mathbf{W} only contains non-negative elements, a sufficient condition for $\rho(\mathbf{W}) < 1$ is that the matrix is sub-stochastic, which is always achieved for $\theta > 0$ and at least one $c_{ij} > 0$ when $a_{ij} > 0$ (see Equation (5)), which is assumed for now. The matrix \mathbf{Z} is called the **fundamental matrix** and z_{ij} is the element i, j of \mathbf{Z} .

Hence, following Equations (6-7), we finally obtain, for the numerator of Equation (4),

$$\sum_{\varphi \in \mathcal{P}_{ij}} \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)] = z_{ij}. \quad (8)$$

On the other hand, for the denominator of Equation (4), we have

$$\sum_{i,j=1}^n \sum_{\varphi \in \mathcal{P}_{ij}} \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)] = \sum_{i,j=1}^n z_{ij} \triangleq \mathcal{Z}, \quad (9)$$

where \mathcal{Z} is called the **partition function**.

Therefore, from Equation (4), the probability of picking a path starting in i end ending in j in our bag of paths model is

$$\mathbf{P}(s = i, e = j) = \frac{z_{ij}}{\mathcal{Z}}, \quad \text{with } \mathbf{Z} = (\mathbf{I} - \mathbf{W})^{-1}. \quad (10)$$

Notice that $\mathbf{P}(s = i, e = j)$ is not symmetric¹. We now turn to a variant of the bag-of-paths, the **bag-of-hitting-paths**.

C. The bag-of-hitting-paths framework

The idea behind the bag-of-hitting-paths model is the same as the bag-of-paths model but the set of paths is now restricted to paths in which the ending node does not appear more than once. In other words, no intermediate node on the path is allowed to be the ending node j (node j is made *absorbing*) and the motifs are now the hitting paths. Hitting paths will play

¹For a symmetric variant in the case of undirected graphs [1], we can consider the probability of picking either $i \rightsquigarrow j$ or $j \rightsquigarrow i$, which is $\mathbf{P}(s = i, e = j) + \mathbf{P}(s = j, e = i)$.

an important role in the derivation of the BoP betweenness. Each non-hitting path $\wp_{ij} \in \mathcal{P}_{ij}$ can be split uniquely into two sub-paths, before hitting node j for the first time, $\wp_{ij}^h \in \mathcal{P}_{ij}^h$ (the set of all hitting paths), and after hitting node j , $\wp_{jj} \in \mathcal{P}_{jj}$ (see [1] for details). Notice the usage of the superscript h to refer to *hitting* paths. The composition of the two sub-paths is a valid path, where $\wp_{ij}^h \circ \wp_{jj} \in \mathcal{P}_{ij}$ is the concatenation of the two paths.

In the case of a bag containing hitting paths, the probability of picking a path $i \rightsquigarrow j$ is defined in a similar way as for non-hitting paths (Equation (4)),

$$P_h(s = i, e = j) = \frac{\sum_{\wp \in \mathcal{P}_{ij}^h} \tilde{\pi}^{\text{ref}}(\wp) \exp[-\theta \tilde{c}(\wp)]}{\sum_{\wp' \in \mathcal{P}^h} \tilde{\pi}^{\text{ref}}(\wp') \exp[-\theta \tilde{c}(\wp')]} \quad (11)$$

and is called the **bag-of-hitting-paths probability** distribution.

Now, since $\tilde{c}(\wp_{ij}) = \tilde{c}(\wp_{ij}^h) + \tilde{c}(\wp_{jj})$ and $\tilde{\pi}^{\text{ref}}(\wp_{ij}) = \tilde{\pi}^{\text{ref}}(\wp_{ij}^h) \tilde{\pi}^{\text{ref}}(\wp_{jj})$ for any $\wp_{ij} = \wp_{ij}^h \circ \wp_{jj}$, we easily obtain

$$\begin{aligned} z_{ij} &= \sum_{\wp_{ij} \in \mathcal{P}_{ij}} \tilde{\pi}^{\text{ref}}(\wp_{ij}) \exp[-\theta \tilde{c}(\wp_{ij})] \\ &= \sum_{\wp_{ij}^h \in \mathcal{P}_{ij}^h} \sum_{\wp_{jj} \in \mathcal{P}_{jj}} \tilde{\pi}^{\text{ref}}(\wp_{ij}^h) \exp[-\theta \tilde{c}(\wp_{ij}^h)] \\ &\quad \times \tilde{\pi}^{\text{ref}}(\wp_{jj}) \exp[-\theta \tilde{c}(\wp_{jj})] \\ &= z_{ij}^h z_{jj}. \end{aligned} \quad (12)$$

From which we deduce $\sum_{\wp \in \mathcal{P}_{ij}^h} \tilde{\pi}^{\text{ref}}(\wp) \exp[-\theta \tilde{c}(\wp)] = z_{ij}^h / z_{jj}$. Now, by analogy with Equation (10), but for hitting paths (Equation (11)),

$$P_h(s = i, e = j) = \frac{z_{ij}^h}{\sum_{i'=1}^n \sum_{j'=1}^n z_{i'j'}^h} = \frac{z_{ij}^h}{Z_h} \quad (13)$$

and the partition function for the bag-of-hitting-paths is therefore

$$Z_h = \sum_{i,j=1}^n z_{ij}^h = \sum_{i,j=1}^n \frac{z_{ij}}{z_{jj}}. \quad (14)$$

Finally, let us just mention that another derivation is available in [1], where it is further shown that z_{ij}^h can be interpreted as either

- The expected reward endorsed by an agent (the reward along a path \wp being defined as $\exp[-\theta \tilde{c}(\wp)]$) when traveling from i to j along all possible paths $\wp \in \mathcal{P}_{ij}^h$ with probability $\tilde{\pi}^{\text{ref}}(\wp)$.
- The expected number of passages through node j for an evaporating, or killed, random walker starting in node i and walking according to the sub-stochastic transition probabilities $p_{ij}^{\text{ref}} \exp[-\theta c_{ij}]$.

Before deriving the BoP betweenness measure, let us consider some related work.

III. RELATED WORK

Graph-based semi-supervised classification has been the subject of intensive research in recent years and a wide range of approaches has been developed in order to tackle the problem [6], [11], [12]: Random-walk-based methods [13], [14], spectral methods [15], [16], regularization frameworks [4], [17]–[19], transductive and spectral SVM [20], to name a few. We will compare our method (the BoP) to some of those techniques, namely,

- 1) A simple alignment with the regularized laplacian kernel (RL) based on a sum of similarities, $\mathbf{K} \mathbf{y}_c$, where $\mathbf{K} = (\mathbf{I} + \lambda \mathbf{L})^{-1}$, $\mathbf{L} = \mathbf{D} - \mathbf{A}$ is the laplacian matrix, \mathbf{I} is the identity matrix, \mathbf{D} is the generalized outdegree matrix, and \mathbf{A} is the adjacency matrix of G [18], [21]. The similarity is computed for each class c in turn. Then, each node is assigned to the class showing the largest sum of similarities. The (scalar) parameter $\lambda > 0$ is the regularization parameter [8], [22].
- 2) A simple alignment with the regularized normalized laplacian kernel (RNL) based on a sum of similarities, $\mathbf{K} \mathbf{y}_c$, where $\mathbf{K} = (\mathbf{I} + \lambda \tilde{\mathbf{L}})^{-1}$, and $\tilde{\mathbf{L}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$ is the normalized laplacian matrix [4], [23]. The assignment to the classes is the same as for the previous method. The regularized normalized laplacian approach seems less sensitive to the priors of the different classes than the un-normalized regularized laplacian approach (RL) [23].
- 3) A simple alignment with the regularized commute time kernel (RCT) based on a sum of similarities, $\mathbf{K} \mathbf{y}_c$, with $\mathbf{K} = (\mathbf{D} - \alpha \mathbf{A})^{-1}$ [4], [22]. The assignment to the classes is the same as for previous methods. Element i, j of this kernel can be interpreted as the discounted cumulated probability of visiting node j when starting from node i . The (scalar) parameter $\alpha \in]0, 1[$ corresponds to an evaporating or killed random walk where the random walker has a $(1 - \alpha)$ probability of disappearing at each step. This method provided the best results in a recent comparative study on semi-supervised classification [22].
- 4) The harmonic function (HF) approach [5], [11], is closely related to the regularization framework of RL and RNL. Furthermore, it is equivalent and provides the same results as the label propagation algorithm [12] and the wvRN (or pRN) classifier used by the Netkit software as a baseline [24], but Netkit is a more general-purpose toolbox able to tackle more complex situations [25]. As those three algorithms give the same results, we only report HF which appeared first in the literature and is fastest. It is based on a structural contiguity measure that smoothes the predicted values and leads to a model having interesting interpretations in terms of electrical potential and absorbing probabilities in a Markov chain.
- 5) The random walk with restart (RWW) classifier [3], [26] relies on random walks performed on the weighted graph seen as a Markov chain. More precisely, a group betweenness measure is derived for each class, based on the stationary distribution of a random walk restarting

from the labeled nodes belonging to a class of interest. Each unlabeled node is then assigned to the class showing maximal betweenness. In this version [22], the random walker has a probability $(1 - \alpha)$ to be teleported – with a uniform probability – to a node belonging to the class of interest c .

- 6) The discriminative random walks approach (\mathcal{D} -walk or DW1; see [2]) also relies on random walks performed on the weighted graph. As for the RWWR, a group betweenness measure, based on passage times during random walks, is derived for each class. More precisely, a \mathcal{D} -walk is a random walk starting in a labeled node and ending when any node having the same label (possibly the starting node itself) is reached for the first time. During this random walk, the number of visits to any unlabeled node is recorded and corresponds to a group betweenness measure. As for the previous method, each unlabeled node is then assigned to the class showing maximal betweenness.
- 7) A modified version of the \mathcal{D} -walk (or DW2). The only difference is that all elements of the transition matrix \mathbf{P}^{ref} (since the random walks is seen as a Markov chain) are multiplied by $\alpha \in]0, 1]$ so that α can be seen as a probability of continuing the random walk at each time step (and so $(1 - \alpha) \in [0, 1[$ is the probability of stopping the random walk at each step. This defines a killed random walk since $\alpha\mathbf{P}^{\text{ref}}$ is now sub-stochastic.

All these methods will be compared to the bag-of-paths (BoP) developed in the next sections. Notice that the random walker of the random-walk-based methods usually follows too long – and thus irrelevant – paths into account: popular entries are therefore intrinsically favored [27], [28]. The bag-of-paths approach tackles this issue by putting a negative exponential term in (5) and part of its success can be imputed to this fact.

Some authors also considered bounded (or truncated) walks [8], [29], [30] and obtained promising results on large graphs. This approach could also be considered in our framework in order to tackle large networks; this will be investigated in further work.

Tong et al. suggested a method avoiding to take the inverse of an $n \times n$ matrix for computing the random walk with restart measure [26]. They reduce the computing time by partitioning the input graph into smaller communities. Then, a sparse approximate of the random walk with restart is obtained by applying a low rank approximation. This approach suffers from the fact that it adds a hyperparameter k (the number of communities) that depends on the network and is still untractable for large graphs with millions of nodes. On the other hand, the computing time is reduced by this same factor k . This is another track to investigate in further work.

Herbster et al. [31] proposed a technique for fast label prediction on graphs through the approximation of the graph with either a minimum spanning tree or a shortest path tree. Once the tree has been extracted, the pseudo inverse of the laplacian matrix can be computed efficiently. The fast computation of the pseudo inverse enables to address prediction problems on large graphs. Finally, Tang and Liu have investigated relational learning via latent social dimensions [32]–[34]. They

proposed to extract latent social dimensions based on network information (such as Facebook, Twitter,...) first, then they used these as features for discriminative learning (via an SVM, for example [32]). Their approach tackles very large networks and provides promising results, especially when only a few labeled data are available.

A lot of research has also been done on collective classification of nodes in networks (see [35] for an introduction). Collective classification also uses the graph topology and a proportion of labeled nodes to classify unlabeled nodes using the same assumption as our proposed technique (i.e. local consistency or homophily).

We also experimented a group betweenness using Freeman’s, i.e. the shortest path, betweenness [36] and a modified version of Newman’s betweenness [37]. For this latter one, the transition probabilities were set to \mathbf{P}^{ref} , and the ending node of the walk was forced to be absorbing. Then, the expected number of visits to each node was recorded and cumulated for each input-output path. However, our BoP group betweenness outperformed these two other class betweenness measures and consequently, results are not reported in this paper.

IV. THE BAG-OF-PATHS BETWEENNESSES

In order to define the BoP classifier, we need to introduce the BoP group betweenness centrality. This concept is itself an extension of the BoP betweenness centrality, which will be developed in the next subsection. The BoP betweenness is related to well-known betweenness measures in some sense: if $\theta \rightarrow \infty$ the BoP betweenness tends to be correlated with Freeman’s betweenness [36] (only shortest-paths are considered), while if $\theta \rightarrow 0^+$, the BoP betweenness tends to be correlated with Newman’s betweenness [37] (based on a natural random walk). This section starts with the presentation of the BoP betweenness centrality measure in Subsection IV-A. Then, its extension, the BoP group betweenness centrality, is described in Subsection IV-B.

A. The bag-of-paths betweenness centrality

The BoP betweenness will measure to which extent a node j is likely to lie in-between other pairs of nodes (i, k) , and therefore is an important intermediary between nodes. In short, the **bag-of-paths betweenness measure** is defined as

$$\text{bet}_j = \sum_{i=1}^n \sum_{k=1}^n \mathbf{P}(\text{int} = j | s = i, e = k; i \neq j \neq k \neq i) \quad (15)$$

which corresponds to the a posteriori probability of finding intermediate node j on a path $i \rightsquigarrow j$ drawn from the bag of paths, cumulated over all source-destination pairs (i, k) .

For computing this quantity from the bag-of-paths framework, we first have to calculate the probability $\mathbf{P}(s = i, \text{int} = j, e = k; i \neq j \neq k \neq i)$ – or \mathbf{P}_{ijk} in short – that such paths

visit an *intermediate* node $int = j$ with $i \neq j \neq k \neq i$. Indeed, by using Equations (3) and (4),

$$\begin{aligned} P_{ijk} &= \sum_{\varphi \in \mathcal{P}_{ik}} \delta(j \in \varphi) P(\varphi) \\ &= \frac{\sum_{\varphi \in \mathcal{P}_{ik}} \delta(j \in \varphi) \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)]}{\sum_{\varphi' \in \mathcal{P}} \tilde{\pi}^{\text{ref}}(\varphi') \exp[-\theta \tilde{c}(\varphi')]} \\ &= \frac{1}{Z} \sum_{\varphi \in \mathcal{P}_{ik}} \delta(j \in \varphi) \tilde{\pi}^{\text{ref}}(\varphi) \exp[-\theta \tilde{c}(\varphi)], \end{aligned} \quad (16)$$

where $\delta(j \in \varphi) = 1$ when node j is visited on path φ , and 0 otherwise.

We will now use the fact that each path φ_{ik} between i and k passing through j can be decomposed uniquely into a *hitting* sub-path φ_{ij} from i to j and a regular sub-path φ_{jk} from j to k (see Subsection II-C). The sub-path φ_{ij} is found by following path φ_{ik} until reaching j for the first time. Therefore, for $i \neq j \neq k \neq i$,

$$\begin{aligned} P_{ijk} &= \frac{1}{Z} \sum_{\varphi_{ij}^h \in \mathcal{P}_{ij}^h} \sum_{\varphi_{jk} \in \mathcal{P}_{jk}} \tilde{\pi}^{\text{ref}}(\varphi_{ij}^h) \tilde{\pi}^{\text{ref}}(\varphi_{jk}) \\ &\quad \times \exp[-\theta \tilde{c}(\varphi_{ij}^h)] \exp[-\theta \tilde{c}(\varphi_{jk})]. \end{aligned} \quad (17)$$

This equation can be reordered to get, for $i \neq j \neq k \neq i$:

$$\begin{aligned} P_{ijk} &= \frac{1}{Z} \left[\sum_{\varphi_{ij}^h \in \mathcal{P}_{ij}^h} \tilde{\pi}^{\text{ref}}(\varphi_{ij}^h) \exp[-\theta \tilde{c}(\varphi_{ij}^h)] \right] \\ &\quad \times \left[\sum_{\varphi_{jk} \in \mathcal{P}_{jk}} \tilde{\pi}^{\text{ref}}(\varphi_{jk}) \exp[-\theta \tilde{c}(\varphi_{jk})] \right]. \end{aligned} \quad (18)$$

Then, after multiplying by Z_h/Z_h , we obtain

$$\begin{aligned} P_{ijk} &= \frac{1}{Z} z_{ij}^h z_{jk} = Z_h \frac{z_{ij}^h}{Z_h} \frac{z_{jk}}{Z} \\ &= Z_h P_h(s = i, e = j) P(s = j, e = k), \end{aligned} \quad (19)$$

with $i \neq j \neq k \neq i$ and where we used Equations (12) and (13).

Finally, recalling Equations (10), (13),

$$\begin{aligned} P_{ijk} &= \frac{\left(\frac{z_{ij}}{z_{jj}} \right) (z_{jk})}{Z} \\ &= \frac{1}{Z} \frac{z_{ij} z_{jk}}{z_{jj}}, \text{ with } i \neq j \neq k \neq i. \\ &= \frac{1}{Z} \frac{z_{ij} z_{jk}}{z_{jj}} \delta(i \neq j \neq k \neq i) \end{aligned} \quad (20)$$

Now, using the Bayes's rule, the *a posteriori* probability of finding intermediate node j on a path starting in i and ending in k is

$$\begin{aligned} P(int = j | s = i, e = k; i \neq j \neq k \neq i) \\ &= \frac{P(s = i, int = j, e = k; i \neq j \neq k \neq i)}{\sum_{j'=1}^n P(s = i, int = j', e = k; i \neq j' \neq k \neq i)} \end{aligned}$$

Using Equation (20), if we assume that node k can be reached from node i , this leads to

$$\begin{aligned} P(int = j | s = i, e = k; i \neq j \neq k \neq i) \\ &= \frac{\left(\frac{z_{ij} z_{jk}}{z_{jj}} \right)}{\sum_{\substack{j'=1 \\ j' \notin \{i,k\}}}^n \left(\frac{z_{ij'} z_{j'k}}{z_{j'j'}} \right)} \delta(i \neq j \neq k \neq i). \end{aligned} \quad (21)$$

Based on this *a posteriori* probability distribution, the *bag-of-paths betweenness* of node j is defined as the sum of the *a posteriori* probabilities of visiting j for all possible starting-ending pairs (i, k) :

$$\begin{aligned} \text{bet}_j &= \sum_{i=1}^n \sum_{k=1}^n P(int = j | s = i, e = k; i \neq j \neq k \neq i) \\ &= \frac{1}{z_{jj}} \sum_{\substack{i=1 \\ i \neq j}}^n \sum_{\substack{k=1 \\ k \notin \{i,j\}}}^n \frac{z_{ij} z_{jk}}{\sum_{\substack{j'=1 \\ j' \notin \{i,k\}}}^n \left(\frac{z_{ij'} z_{j'k}}{z_{j'j'}} \right)}, \end{aligned} \quad (23)$$

which allows to compute the betweenness from the fundamental matrix \mathbf{Z} (Equation (7)).

Let us now derive the matrix formula providing the betweenness vector bet , containing the betweennesses for each node. First of all, the normalization factor appearing in the denominator of Equation (23), denoted here by n_{ik} , is computed,

$$n_{ik} = \sum_{j'=1}^n (1 - \delta_{ij'}) (1 - \delta_{j'k}) (z_{ij'} z_{j'k}) / z_{j'j'}, \quad (24)$$

which can be re-written as

$$n_{ik} = \sum_{j'=1}^n \{(1 - \delta_{ij'}) z_{ij'}\} \{1/z_{j'j'}\} \{(1 - \delta_{j'k}) z_{j'k}\}. \quad (25)$$

Therefore, the matrix containing the normalization factors n_{ik} is

$$\mathbf{N} = (\mathbf{Z} - \text{Diag}(\mathbf{Z})) (\text{Diag}(\mathbf{Z}))^{-1} (\mathbf{Z} - \text{Diag}(\mathbf{Z})), \quad (26)$$

where for a given matrix \mathbf{M} , $\text{diag}(\mathbf{M})$ is a column vector containing the diagonal of \mathbf{M} and $\text{Diag}(\mathbf{M})$ is a diagonal matrix containing the diagonal of \mathbf{M} .

Moreover, the inner term appearing in Equation (23) can be rewritten as

$$\begin{aligned} \sum_{i=1}^n \sum_{k=1}^n \delta(i \neq j \neq k \neq i) z_{ij} (1/n_{ik}) z_{jk} \\ &= \sum_{i=1}^n \sum_{k=1}^n \{(1 - \delta_{ji}) z_{ji}^t\} \{(1 - \delta_{ik})(1/n_{ik})\} \{(1 - \delta_{kj}) z_{kj}^t\}, \end{aligned} \quad (27)$$

where z_{ij}^t is the element i, j of matrix \mathbf{Z}^T (transpose of \mathbf{Z}). In matrix form, bet (see Equation (21)) is therefore equal to

$$\begin{aligned} \text{bet} &= (\text{Diag}(\mathbf{Z}))^{-1} \text{diag}[(\mathbf{Z}^T - \text{Diag}(\mathbf{Z})) \\ &\quad \times (\mathbf{N}^\dagger - \text{Diag}(\mathbf{N}^\dagger)) (\mathbf{Z}^T - \text{Diag}(\mathbf{Z}))], \end{aligned} \quad (28)$$

with matrix \mathbf{N}^\dagger containing elements $n_{ik}^\dagger = 1/n_{ik}$ (element-wise reciprocal).

B. The bag-of-paths group betweenness centrality

Let us now generalize the bag-of-paths betweenness to a group betweenness measure. Quite naturally, the **bag-of-paths group betweenness** of node j will be defined as

$$\text{gbet}_j(\mathcal{C}_i, \mathcal{C}_k) = \text{P}(int = j | s \in \mathcal{C}_i, e \in \mathcal{C}_k; s \neq int \neq e \neq s) \quad (29)$$

and can be interpreted as the extent to which the node j lies in-between the two sets of nodes \mathcal{C}_i and \mathcal{C}_k . It is assumed that the sets $\{\mathcal{C}_i\}_{i=1}^m$ are disjoint. Using Bayes' law provides

$$\begin{aligned} & \text{P}(int = j | s \in \mathcal{C}_i, e \in \mathcal{C}_k; s \neq int \neq e \neq s) \\ &= \frac{\text{P}(s \in \mathcal{C}_i, int = j, e \in \mathcal{C}_k; s \neq int \neq e \neq s)}{\text{P}(s \in \mathcal{C}_i, e \in \mathcal{C}_k; s \neq int \neq e \neq s)} \\ &= \frac{\sum_{i' \in \mathcal{C}_i} \sum_{k' \in \mathcal{C}_k} \text{P}(s = i', int = j, e = k'; s \neq int \neq e \neq s)}{\sum_{j'=1}^n \sum_{i' \in \mathcal{C}_i} \sum_{k' \in \mathcal{C}_k} \text{P}(s = i', int = j', e = k'; s \neq int \neq e \neq s)} \end{aligned} \quad (30)$$

Substituting (20) for the joint probabilities in Equation (30) allows to compute the group betweenness measure in terms of the elements of the fundamental matrix \mathbf{Z} :

$$\begin{aligned} & \text{gbet}_j(\mathcal{C}_i, \mathcal{C}_k) \\ &= \frac{\sum_{i' \in \mathcal{C}_i} \sum_{k' \in \mathcal{C}_k} \delta(i' \neq j \neq k' \neq s) \frac{z_{i'j} z_{jk'}}{z_{jj}}}{\sum_{j'=1}^n \sum_{i' \in \mathcal{C}_i} \sum_{k' \in \mathcal{C}_k} \delta(i' \neq j' \neq k' \neq s) \frac{z_{i'j'} z_{j'k'}}{z_{j'j'}}}, \end{aligned} \quad (31)$$

where the denominator is simply a normalization factor ensuring that the probability distribution sums to one. It is therefore sufficient to compute the numerator only and then normalize the resulting quantity.

Let us put this expression in matrix form. As before, we denote element i, j of matrix \mathbf{Z}^T as z_{ij}^T . It is also assumed that node i' and k' belong to different groups, $\mathcal{C}_i \neq \mathcal{C}_k$, so that i and k are necessarily different (classes are disjoint). The numerator in Equation (31) is

$$\begin{aligned} & \text{num}(\text{gbet}_j(\mathcal{C}_i, \mathcal{C}_k)) \\ &= \frac{1}{z_{jj}} \sum_{i' \in \mathcal{C}_i} \sum_{k' \in \mathcal{C}_k} (1 - \delta_{ji'}) (1 - \delta_{jk'}) z_{i'j} z_{jk'} \\ &= \frac{1}{z_{jj}} \left(\sum_{i' \in \mathcal{C}_i} (1 - \delta_{ji'}) z_{ji'}^T \right) \left(\sum_{k' \in \mathcal{C}_k} (1 - \delta_{jk'}) z_{jk'} \right). \end{aligned} \quad (32)$$

If y_i^c is a binary indicator indicating if node i belongs to the class c (as described in Section II-A), the numerator can be rewritten as

$$\begin{aligned} & \text{num}(\text{gbet}_j(\mathcal{C}_i, \mathcal{C}_k)) \\ &= \frac{1}{z_{jj}} \left(\sum_{i'=1}^n (1 - \delta_{ji'}) z_{ji'}^T y_{i'}^i \right) \left(\sum_{k'=1}^n (1 - \delta_{jk'}) z_{jk'} y_{k'}^k \right). \end{aligned} \quad (33)$$

Consequently, in matrix form, the group betweenness vector reads

$$\begin{cases} \text{gbet}(\mathcal{C}_i, \mathcal{C}_k) \leftarrow (\mathbf{Diag}(\mathbf{Z}))^{-1} ((\mathbf{Z}_0^T \mathbf{y}^i) \circ (\mathbf{Z}_0 \mathbf{y}^k)) \\ \quad \text{with } \mathbf{Z}_0 = \mathbf{Z} - \mathbf{Diag}(\mathbf{Z}), \\ \text{gbet}(\mathcal{C}_i, \mathcal{C}_k) \leftarrow \frac{\text{gbet}(\mathcal{C}_i, \mathcal{C}_k)}{\|\text{gbet}(\mathcal{C}_i, \mathcal{C}_k)\|_1} \text{ (normalization)} \end{cases} \quad (34)$$

where we assume $i \neq k$. In this equation, the vector $\text{gbet}(\mathcal{C}_i, \mathcal{C}_k)$ must be normalized by dividing it by its L_1 norm. Notice that $\mathbf{Z}_0 = \mathbf{Z} - \mathbf{Diag}(\mathbf{Z})$ is simply the fundamental matrix whose diagonal is set to zero.

V. SEMI-SUPERVISED CLASSIFICATION THROUGH THE BAG-OF-PATHS GROUP BETWEENNESS

In this section, the bag-of-paths model, and more precisely the bag-of-paths group betweenness measure, will be used for *classification purposes*. Notice, however, that in the derivation of the group betweenness measure (see Equation (34)), it was assumed that the starting and ending classes are different ($\mathcal{C}_i \neq \mathcal{C}_k$). We will now recompute this quantity when starting and ending in the same class c , i.e. calculating $\text{gbet}_j(\mathcal{C}_c, \mathcal{C}_c)$, which provides a **within-class betweenness**. Indeed, this quantity measures the extent to which nodes of G are in-between – and therefore in the neighborhood of – the nodes of class c .

A within-class betweenness is thus computed for each class c and each node will define assigned to the class showing the highest betweenness. This will be our simple classification rule based on the within-class betweenness. The main hypothesis underlying this classification technique is that a node is likely to belong to the same class as its “neighboring nodes”. This is usually called the local consistency assumption (also called smoothness, homophily or cluster assumption [5], [12], [38]).

The same reasoning as for deriving Equation (34) is applied in order to compute the numerator of (31) in this new case. We start with Equation (31), considering now the same starting and ending class c but multiplying the term inside the double sum by $(1 - \delta_{i'k'})$. This new term will ensure that the starting node is different from the ending node (this was always the case with different starting and ending classes, but now this must be forced). From Equation (32), this can be rewritten as

$$\begin{aligned} \text{num}(\text{gbet}_j(\mathcal{C}_c, \mathcal{C}_c)) &= \frac{1}{z_{jj}} \sum_{i', k' \in \mathcal{C}_c} (1 - \delta_{ji'}) (1 - \delta_{i'k'}) \\ &\quad \times (1 - \delta_{jk'}) z_{i'j} z_{jk'}. \end{aligned} \quad (35)$$

$\text{num}(\text{gbet}_j(\mathcal{C}_c, \mathcal{C}_c))$ is thus the same as $\text{num}(\text{gbet}_j(\mathcal{C}_i, \mathcal{C}_k))$ with $\mathcal{C}_i = \mathcal{C}_c$ of Equation (32) and $\mathcal{C}_k = \mathcal{C}_c$, plus an extra term:

$$\begin{aligned} \text{num}(\text{gbet}_j(\mathcal{C}_c, \mathcal{C}_c)) &= \\ &= \frac{1}{z_{jj}} \sum_{i' \in \mathcal{C}_c} \sum_{k' \in \mathcal{C}_c} (1 - \delta_{ji'}) (1 - \delta_{jk'}) z_{i'j} z_{jk'} \\ &\quad - \frac{1}{z_{jj}} \sum_{i' \in \mathcal{C}_c} \sum_{k' \in \mathcal{C}_c} (1 - \delta_{ji'}) \delta_{i'k'} (1 - \delta_{jk'}) z_{i'j} z_{jk'}. \end{aligned} \quad (36)$$

TABLE I
CLASS DISTRIBUTION OF THE NINE *NewsGroups* DATA SETS. NG 1-3 CONTAIN TWO CLASSES, NG 4-6 CONTAIN THREE CLASSES AND NG 7-9 CONTAIN FIVE CLASSES.

Class	NG1	NG2	NG3	NG4	NG5	NG6	NG7	NG8	NG9
1	200	198	200	200	200	197	200	200	200
2	200	200	199	200	198	200	200	200	200
3				200	200	198	200	198	197
4							200	200	200
5							198	200	200
Total	400	398	399	600	598	595	998	998	997

It is easy to show that with this additional term, the matrix equation for $\text{num}(\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c))$ (Equation (34)) becomes

$$\text{num}(\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)) = (\mathbf{Diag}(\mathbf{Z}))^{-1} [(\mathbf{Z}_0^T \mathbf{y}^c) \circ (\mathbf{Z}_0 \mathbf{y}^c) - (\mathbf{Z}_0^T \circ \mathbf{Z}_0) \mathbf{y}^c]. \quad (37)$$

Once again, this is the same result as in Equation (34) with one more term that ensure that the starting node is different from the ending node. After having computed this equation, the numerator must be normalized in order to obtain $\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)$ (as for Equation (34)).

Finally, if we want to classify a node, $\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)$ is computed for each class c in turn and then, for each node, the class label showing the maximal betweenness is chosen,

$$\hat{\ell} = \arg \max_{c \in \mathcal{L}} (\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)), \text{ with}$$

$$\begin{cases} \mathbf{D}_z = \mathbf{Diag}(\mathbf{Z}); \mathbf{Z}_0 = \mathbf{Z} - \mathbf{D}_z \text{ (set diagonal to 0)} \\ \mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c) \leftarrow \mathbf{D}_z^{-1} [(\mathbf{Z}_0^T \mathbf{y}^c) \circ (\mathbf{Z}_0 \mathbf{y}^c) - (\mathbf{Z}_0^T \circ \mathbf{Z}_0) \mathbf{y}^c] \\ \mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c) \leftarrow \frac{\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)}{\|\mathbf{gbet}(\mathcal{C}_c, \mathcal{C}_c)\|_1} \text{ (normalization)} \end{cases} \quad (38)$$

where \mathcal{L} is the set of class labels. The pseudo-code for the BoP classifier can be found in Algorithm 1. Of course, once computed, the group betweenness is only used to classify the unlabeled nodes.

VI. EXPERIMENTAL COMPARISONS

In this section, the bag-of-paths group betweenness approach for semi-supervised classification (referred to as the BoP classifier for simplicity) will be compared to other semi-supervised classification techniques on multiple data sets. The different classifiers to which the BoP classifier will be compared were already introduced in Section III and are recalled in Table IV.

TABLE II
CLASS DISTRIBUTION OF THE *IMDb-proco* DATA SET.

Class	IMDb
High-revenue	572
Low-revenue	597
Total	1169

The goal of the experiments of this section is to classify unlabeled nodes in partially labeled graphs and to compare the different methods in terms of classification accuracy. This

Algorithm 1 Classification through the bag-of-paths group betweenness algorithm.

Input:

- A weighted directed graph G containing n nodes, represented by its $n \times n$ adjacency matrix \mathbf{A} , containing affinities.
- The $n \times n$ transition cost matrix \mathbf{C} associated to G .
- m binary indicator vectors \mathbf{y}_c containing as entries 1 for nodes belonging to the class \mathcal{C}_c , and 0 otherwise. Classes are mutually exclusive.
- The inverse temperature parameter θ .

Output:

- The $n \times 1$ vector $\hat{\ell}$ containing the predicted class labels of each node.

- 1: $\mathbf{D} \leftarrow \mathbf{Diag}(\mathbf{Ae})$ {the row-normalization matrix}
- 2: $\mathbf{P}^{\text{ref}} \leftarrow \mathbf{D}^{-1} \mathbf{A}$ {the reference transition probabilities matrix}
- 3: $\mathbf{W} \leftarrow \mathbf{P}^{\text{ref}} \circ \exp[-\theta \mathbf{C}]$ {elementwise exponential and multiplication \circ }
- 4: $\mathbf{Z} \leftarrow (\mathbf{I} - \mathbf{W})^{-1}$ {the fundamental matrix}
- 5: $\mathbf{Z}_0 \leftarrow \mathbf{Z} - \mathbf{Diag}(\mathbf{Z})$ {set diagonal to zero}
- 6: $\mathbf{D}_z \leftarrow \mathbf{Diag}(\mathbf{Z})$
- 7: **for** $c = 1$ to m **do**
- 8: $\hat{\mathbf{y}}_c^* \leftarrow \mathbf{D}_z^{-1} [(\mathbf{Z}_0^T \mathbf{y}^c) \circ (\mathbf{Z}_0 \mathbf{y}^c) - (\mathbf{Z}_0^T \circ \mathbf{Z}_0) \mathbf{y}^c]$ {compute the group betweenness for class c ; \circ is the elementwise multiplication (Hadamard product)}
- 9: $\hat{\mathbf{y}}_c^* \leftarrow \frac{\hat{\mathbf{y}}_c^*}{\|\hat{\mathbf{y}}_c^*\|_1}$ {normalize the betweenness scores}
- 10: **end for**
- 11: $\hat{\ell} \leftarrow \arg \max_{c \in \mathcal{L}} (\hat{\mathbf{y}}_c^*)$ {each node is assigned to the class showing the largest class betweenness}
- 12: **return** $\hat{\ell}$

TABLE III
CLASS DISTRIBUTION OF THE FOUR *WebKB cocite* DATA SETS.

Class	Cornell	Texas	Washington	Wisconsin
Course	54	51	170	83
Department	25	36	20	37
Faculty	62	50	44	37
Project	54	28	39	25
Staff	6	6	10	11
Student	145	163	151	155
Total	346	334	434	348
Majority class (%)	41.9	48.8	39.2	44.5

comparison is performed on medium-size networks only since kernel approaches are difficult to compute on large networks. The computational tractability of the methods used in this experimental section will also be analyzed.

This section is organized as follows. First, the data sets used for the semi-supervised classification will be described

TABLE IV

THE EIGHT CLASSIFIERS, THE VALUE RANGE TESTED FOR TUNING THEIR PARAMETERS AND THE MOST SELECTED VALUES. MODE1 IS THE MOST SELECTED VALUE, MODE2 IS THE SECOND MOST SELECTED VALUE AND MODE3 IS THE THIRD MOST SELECTED VALUE. NOTICE THAT DW2 WITH $\alpha = 1.0$ IS THE SAME AS DW1.

Classifier name	Acronym	Parameter	Tested values	Mode1	Mode2	Mode3
Regularized laplacian kernel	RL	$\lambda > 0$	$10^{-6}, 10^{-5}, \dots, 10^6$	10^{-9} (12.3%)	10^{-1} (11.7%)	10^{-2} (11.5%)
Regularized normalised laplacian kernel	RNL	$\lambda > 0$	$10^{-6}, 10^{-5}, \dots, 10^6$	10^{-1} (42.1%)	10^{-2} (13.9%)	10^{-3} (09.3%)
Label Propagation	LP	none	/	/	/	/
Regularized commute-time kernel	RCT	$\alpha \in [0, 1]$	0.1, 0.2, ..., 1	0.9(28.0%)	0.8(16.2%)	0.7(12.2%)
Harmonic function	HF	none	/	/	/	/
Random walk with restart	RWWR	$\alpha \in [0, 1]$	0.1, 0.2, ..., 1	0.9(45.8%)	0.8(16.8%)	0.7(10.1%)
Discriminative random walks	DW1	none	/	/	/	/
Killed discriminative random walks	DW2	$\alpha \in [0, 1]$	0.1, 0.2, ..., 1	1.0(19.5%)	0.1(11.8%)	0.9(11.2%)
BoP classifier	BoP	$\theta > 0$	$10^{-6}, 10^{-5}, \dots, 10^2$	10^{-4} (28.3%)	10^{-3} (25.9%)	10^{-2} (12.3%)

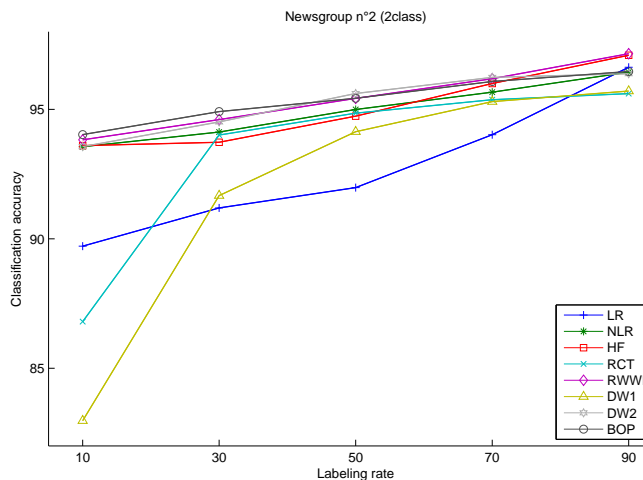


Fig. 1. Classification accuracies in percents, averaged over 20 runs, obtained on partially labeled graphs. Results are reported for the eight methods (RL, RNL, RCT, HF, RWWR, DW1, DW2, BoP) and for five labeling rates (10%, 30%, 50%, 70%, 90%). This graph shows the results obtained on the *NewsGroups* (NG2) data set.

in Subsection VI-A. Second, the experimental methodology is detailed in Subsection VI-B. Third, the results will be discussed in Subsection VI-C. Fourth, the computation time will be investigated in Subsection VI-D. Finally, extreme cases are studied in Subsection VI-E.

A. Datasets

The different classifiers are compared on 14 data sets that have been used previously for semi-supervised classification: nine *NewsGroups* data sets [39], the four universities *WebKB cocite* data sets [24], [17] and the *IMDb prodcoc* data set [24]².

NewsGroups: The *NewsGroups* data set is composed of about 20,000 unstructured documents, taken from 20 discussion groups (newsgroups) of the Usenet diffusion list. 20 Classes (or topics) were originally present in the data set. For our experiments, nine subsets related to different topics are extracted from the original data set, resulting in a total of nine different

²The different data sets used for these comparisons are described in Subsection VI-A. Implementations and datasets are available at <http://www.isys.ucl.ac.be/staff/lebicot/research.htm>.

data sets. The data sets were built by sampling about 200 documents at random in each topic (three samples of two, three and five classes, thus nine samples in total). The repartition is listed in Table I. The extraction process and the procedure used for building the graph are detailed in [40].

WebKB cocite: These data sets consist of sets of web pages gathered from four computer science departments (four data sets, one for each university), with each page manually labeled into one of six categories: course, department, faculty, project, staff, and student [24]. The pages are linked by co-citation (if x links to z and y links to z , then x and y are co-citing z), resulting in an undirected graph. The composition of the data sets is shown in Table III.

IMDb-prodcoc: The collaborative Internet Movie Database (*IMDb*, [24]) has several applications such as making movie recommendations or movie category classification. The classification problem focuses on the prediction of whether the movie is a box-office hit or not. It contains a graph of movies linked together whenever they share the same production company and weight of an edge in the graph is the number of production companies that two movies have in common. The *IMDb-prodcoc* class distribution is shown in Table II.

B. Experimental methodology

The classification accuracy will be reported for several labeling rates (10%, 30%, 50%, 70%, 90%), i.e. proportions of nodes for which the label is known. The labels of remaining nodes are deleted during the modeling phase and are used as test data during the assessment phase. For each considered labeling rate, 20 random node label deletions were performed (20 runs) and performances are averaged on these 20 runs. For each unlabeled node, the various classifiers predict the most suitable category. For each run, a 10-fold nested cross-validation is performed for tuning the parameters of the models. The external folds are obtained by 10 successive rotations of the nodes and the performance of one specific run is the average over these 10 folds. Moreover, for each fold of the external cross-validation, a 10-fold internal cross-validation is performed on the remaining labeled nodes in order to tune the hyperparameters of the classifiers (i.e. parameters α , λ and θ (see Table IV) – methods HF and DW1 do not have any hyperparameter). Thus, for each method and each labeling rate, the mean classification accuracy averaged on the 20 runs will be reported.

C. Results & discussion

Comparative results for each method on the fourteen data sets are reported as follows: the results on the nine *NewsGroups* data sets, on the four *WebKB Cocite* data sets and on the *IMBd-prodco* data set are shown in Table V. The results of the second *NewsGroups* data set are also reported as a plot on Fig. 1 to visualize the typical relation between classification accuracy and labeling rate.

Statistical significance tests for each labeling rate are detailed in Table VI. One-sided *t*-tests were performed to determine whether or not the performance of a method is significantly superior (p-value less than 0.05 on the 20 runs) to another. Table VI can be read as follows. At the bottom of each table, the Win/Tie/Lose frequency summarizes how many times the BoP classifier was significantly better (Win), equivalent (Tie), or significantly worse (Lose) than each other method.

The most selected values of the parameters for each method are reported in Table IV. With each methods, all the tested values were selected at least once. RNL, RCT, RWWR and BoP selected a short range of parameters, which is of course desirable. Some methods, such as RL and DW2 selected a very wide range of parameters, and the different most represented values are not grouped.

Moreover, for each labeling rate, the different classifiers have been ordered according to a Borda score ranking. For each data set, each method is granted with a certain number of points, or rating. This number of points is equal to eight if the classifier is the best classifier (i.e., has the best mean classification accuracy on this data set), seven if the classifier is the second best and so on, so that the worst classifier is granted with only one point. The ratings are then summed across all the considered data sets and the classifiers are sorted by descending total rating. The final ranking, together with the total ratings, are reported in Table VII.

We observe that the BoP classifier always achieved competitive results since it ranges among the top methods on all data sets. The BoP classifier actually tends to be the best algorithm for all labeling rates except for the 90% labeling rate, where it comes second as observed from Table VII and from Table VI. The RWWR method is often second. The RCT kernel also achieves good performance and is the best of the kernel-based classifiers (as suggested in [22]). It is also the best algorithm when the labeling rate is very high (90%).

Notice that RWWR, RCT and DW2 largely outperform the other algorithms (besides BoP). However, it is difficult to figure out which of those three methods is the best, after BoP. It can be noticed that the DW2 version of the \mathcal{D} -walk is more competitive when the labeling rate is low and that it performs much better than the DW1 version, especially for low labeling rates: the Win/Tie/Lose scores for DW2 against DW1 are 7/1/6, 5/2/7, 8/1/5, 13/1/0 and 14/0/0 respectively for 90%, 70%, 50%, 30%, 10% labeling rate.

From the fifth to the eighth position, the ranking is less clear since none of the methods is really better than the other. However, all of these methods (NR and RNL as well as HF and DW1) are significantly worse than BoP, RCT, RWWR and

DW2. Notice also that the performance of DW1 and HF drops significantly when labeling rate decreases.

D. Computation time

The computational tractability of a method is an important consideration to take into account. Table VIII provides a comparison of the running time of all methods. To explore computation time with respect to the number of nodes and the number of classes, artificial graphs with a certain number of classes have been created. For each method, 20 runs on each of the data sets are performed and the running time is recorded for each run. The 20 running times are averaged and results are reported in Table VIII.

We observe that HF is one of the quickest methods, but sadly it is not competitive in terms of accuracy, as reported in Subsection (VI-C). Notice that the two kernel methods, RL and RCT, have more or less the same computation time since the alignment is done once for all the classes. RNL, the last kernel method, is slower than RL, HF and RCT (because of the time-consuming normalisation). After the HF and the kernel methods, the BoP classifier achieves competitive results compared with the remaining classifiers. The time augmentation when the graph size increases is similar for all methods (except for RL and RC for which the augmentation is smaller). The cause is that all those methods require a matrix inversion: the complexity of such an operation is $\mathcal{O}(n^3)$ (where n is the number of nodes) and this is what can be observed from Table VIII (when the number of nodes doubles, the time is more or less multiplied by eight). But the BoP classifier has the same advantage as the kernel methods: its computation time does not increase strongly when the number of classes increases. This comes from the algorithm structure: to contrary to RWWR, DW1 and DW2, the BoP classifier does not require a matrix inversion for each class. Furthermore, the matrix inversions (or linear systems of equations to solve) required for the BoP can be computed as far as the graph (through its adjacency matrix) is known, which is not the case with kernel methods. This is a good property for BoP, since it means that rows 1 to 6 of Algorithm 1 can be pre-computed once for all folds in the cross-validation. Finally, the space complexity is $\mathcal{O}(n^2)$ for all the methods.

E. Extreme cases

In this subsection, two extreme classification cases will be studied. First, what happens if only one or two labeled data points are available? As described in Section V, the BoP classifier requires at least two nodes for computing the BoP group betweenness. We performed a small experiment on the first *NewsGroups* NG1 dataset. The parameters were tuned by a 10-fold cross-validation and 20 runs were averaged. The classification accuracies are reported in Table IX. Only the RCT, RWWR and BoP classifiers remain competitive for this first extreme case.

Secondly, let us consider the case where the classes are imbalanced. The following experiment was designed to study this other extreme case. The classes of the well-known *industry-yh* dataset were merged to get two classes: the

TABLE V

CLASSIFICATION ACCURACIES IN PERCENTS \pm THE STANDARD DEVIATION, AVERAGED OVER 20 RUNS, OBTAINED ON PARTIALLY LABELED GRAPHS. RESULTS ARE REPORTED FOR THE EIGHT METHODS (RL, RNL, LP, RCT, HF, RWWR, DW1, DW2, BoP) AND FOR FIVE LABELING RATES (10%, 30%, 50%, 70%, 90%). THE TABLE SHOWS THE RESULTS FOR ALL THE TESTED DATA SETS. THE STANDARD DEVIATION IS CALCULATED OVER THE 10 FOLDS OF THE EXTERNAL CROSS-VALIDATION OF THE 20 INDEPENDENT RUNS.

l	RL	RNL	RCT	HF	RWWR	DW1	DW2	BoP	
NG1	90%	97.56 \pm 2.28	98.42 \pm 2.03	97.16 \pm 2.46	97.39 \pm 2.55	97.17 \pm 2.42	97.16 \pm 2.71	97.56 \pm 2.50	97.58 \pm 2.08
	70%	94.64 \pm 10.11	97.31 \pm 1.74	96.66 \pm 1.26	97.25 \pm 1.35	96.76 \pm 1.32	96.65 \pm 1.50	96.59 \pm 1.51	97.36 \pm 1.04
	50%	92.87 \pm 12.97	96.77 \pm 1.27	96.53 \pm 1.18	96.80 \pm 1.04	96.70 \pm 1.32	95.73 \pm 1.30	96.41 \pm 1.17	97.27 \pm 0.98
	30%	93.79 \pm 10.52	96.31 \pm 1.35	95.99 \pm 1.30	95.92 \pm 1.06	96.14 \pm 1.10	94.34 \pm 1.24	96.03 \pm 0.94	96.94 \pm 0.70
	10%	87.27 \pm 18.01	95.36 \pm 2.17	96.15 \pm 1.09	88.46 \pm 7.32	96.19 \pm 0.63	88.55 \pm 1.94	95.27 \pm 0.84	96.50 \pm 1.27
NG2	90%	96.63 \pm 2.78	96.44 \pm 2.54	97.09 \pm 2.67	95.61 \pm 3.14	97.15 \pm 2.53	95.71 \pm 3.00	96.38 \pm 2.54	96.46 \pm 2.75
	70%	94.02 \pm 9.97	95.66 \pm 1.86	96.00 \pm 1.63	95.37 \pm 1.65	96.18 \pm 1.46	95.30 \pm 1.74	96.23 \pm 1.51	96.07 \pm 1.81
	50%	91.98 \pm 11.63	94.99 \pm 2.22	94.73 \pm 1.43	94.85 \pm 1.30	95.42 \pm 1.37	94.13 \pm 1.54	95.61 \pm 1.50	95.43 \pm 1.71
	30%	91.19 \pm 11.90	94.12 \pm 3.46	93.73 \pm 1.08	94.01 \pm 1.09	94.60 \pm 1.18	91.67 \pm 1.60	94.51 \pm 1.42	94.91 \pm 1.43
	10%	89.72 \pm 13.38	93.55 \pm 3.25	93.60 \pm 0.80	86.80 \pm 6.11	93.82 \pm 0.63	82.97 \pm 2.11	93.56 \pm 1.29	94.02 \pm 1.45
NG3	90%	96.80 \pm 2.62	96.59 \pm 2.16	98.05 \pm 2.58	96.81 \pm 2.76	98.05 \pm 2.62	96.91 \pm 2.83	96.94 \pm 2.83	98.06 \pm 2.34
	70%	95.66 \pm 6.75	96.63 \pm 0.95	97.70 \pm 1.52	96.59 \pm 1.41	97.77 \pm 1.59	96.89 \pm 1.37	97.19 \pm 1.18	97.77 \pm 1.57
	50%	93.00 \pm 12.04	96.43 \pm 0.94	97.13 \pm 1.08	96.54 \pm 0.92	97.21 \pm 1.24	96.70 \pm 0.98	96.43 \pm 0.90	97.36 \pm 1.26
	30%	93.00 \pm 10.84	95.98 \pm 0.84	96.59 \pm 0.81	95.80 \pm 0.93	96.64 \pm 0.81	95.78 \pm 1.06	96.27 \pm 0.91	96.84 \pm 1.00
	10%	87.73 \pm 16.95	95.29 \pm 0.97	95.45 \pm 1.04	83.44 \pm 9.70	95.21 \pm 1.24	91.81 \pm 1.47	95.46 \pm 0.73	96.05 \pm 0.81
NG4	90%	95.14 \pm 2.18	95.99 \pm 2.01	95.58 \pm 2.07	95.03 \pm 2.68	95.55 \pm 2.05	94.99 \pm 2.74	94.73 \pm 2.17	95.19 \pm 2.27
	70%	89.17 \pm 16.52	94.76 \pm 0.96	94.41 \pm 1.01	94.79 \pm 1.44	94.29 \pm 1.02	95.06 \pm 1.45	94.23 \pm 1.00	94.32 \pm 0.92
	50%	87.59 \pm 17.33	93.60 \pm 0.94	93.27 \pm 0.88	94.38 \pm 1.01	93.27 \pm 0.78	94.35 \pm 1.21	93.82 \pm 0.58	93.85 \pm 0.77
	30%	87.22 \pm 18.23	93.26 \pm 0.88	92.97 \pm 0.82	93.00 \pm 1.02	93.16 \pm 0.90	92.65 \pm 1.13	93.38 \pm 0.58	93.51 \pm 0.75
	10%	84.41 \pm 22.69	91.05 \pm 6.17	94.53 \pm 1.38	74.92 \pm 9.69	95.60 \pm 0.52	85.99 \pm 1.63	94.86 \pm 0.52	95.14 \pm 1.06
NG5	90%	95.35 \pm 2.41	95.65 \pm 2.84	95.94 \pm 2.46	95.25 \pm 2.54	95.97 \pm 2.52	94.97 \pm 2.57	95.84 \pm 2.33	95.80 \pm 2.98
	70%	93.83 \pm 7.58	94.40 \pm 1.47	94.78 \pm 1.13	94.83 \pm 1.36	95.05 \pm 1.07	94.67 \pm 1.48	94.12 \pm 1.55	94.74 \pm 1.22
	50%	91.26 \pm 9.54	92.59 \pm 1.87	93.65 \pm 1.36	93.94 \pm 1.02	94.48 \pm 0.80	93.69 \pm 1.17	93.64 \pm 1.38	94.41 \pm 1.00
	30%	87.89 \pm 14.24	90.41 \pm 1.92	92.41 \pm 1.18	91.51 \pm 1.30	93.83 \pm 0.69	91.65 \pm 1.20	92.84 \pm 0.96	94.21 \pm 1.19
	10%	88.76 \pm 11.94	90.64 \pm 2.22	93.58 \pm 1.01	77.94 \pm 6.73	94.84 \pm 0.65	84.42 \pm 1.89	92.35 \pm 0.60	94.27 \pm 1.07
NG6	90%	93.89 \pm 2.91	92.49 \pm 2.50	96.25 \pm 2.84	94.19 \pm 2.81	96.27 \pm 2.84	93.99 \pm 2.93	95.21 \pm 2.63	96.17 \pm 2.77
	70%	92.43 \pm 1.47	91.05 \pm 1.21	95.93 \pm 1.59	93.07 \pm 1.55	96.07 \pm 1.57	93.86 \pm 1.76	94.00 \pm 1.37	95.96 \pm 1.75
	50%	91.44 \pm 1.00	90.19 \pm 1.76	94.80 \pm 1.10	91.42 \pm 1.34	95.41 \pm 0.71	92.90 \pm 1.39	93.43 \pm 1.27	95.36 \pm 0.77
	30%	89.58 \pm 4.11	88.42 \pm 1.69	93.33 \pm 1.21	88.04 \pm 1.83	94.52 \pm 0.85	90.75 \pm 1.49	92.49 \pm 0.82	94.51 \pm 1.19
	10%	89.19 \pm 10.35	89.20 \pm 3.13	94.18 \pm 0.65	73.74 \pm 5.37	94.93 \pm 0.56	83.66 \pm 1.84	92.82 \pm 0.88	94.21 \pm 1.44
NG7	90%	92.48 \pm 2.31	91.18 \pm 2.90	91.98 \pm 2.52	92.22 \pm 2.46	91.99 \pm 2.51	92.46 \pm 2.52	92.61 \pm 2.25	93.06 \pm 2.02
	70%	91.73 \pm 1.51	91.04 \pm 1.46	91.57 \pm 1.19	91.34 \pm 1.58	91.56 \pm 1.14	91.87 \pm 1.44	91.97 \pm 1.25	92.04 \pm 1.15
	50%	89.97 \pm 8.93	90.40 \pm 1.10	91.42 \pm 0.79	90.33 \pm 1.08	91.61 \pm 0.51	90.69 \pm 1.10	91.88 \pm 0.80	91.61 \pm 0.87
	30%	90.65 \pm 1.33	89.86 \pm 0.71	90.65 \pm 1.16	87.05 \pm 1.32	91.06 \pm 0.68	88.45 \pm 1.09	91.09 \pm 1.12	90.97 \pm 0.75
	10%	85.52 \pm 16.28	88.66 \pm 0.76	88.89 \pm 1.10	68.85 \pm 6.90	91.08 \pm 0.57	80.31 \pm 1.67	90.33 \pm 0.57	90.23 \pm 1.05
NG8	90%	90.34 \pm 2.70	90.16 \pm 2.12	90.81 \pm 2.42	88.73 \pm 2.79	90.80 \pm 2.43	88.73 \pm 2.87	90.54 \pm 2.19	90.48 \pm 2.81
	70%	89.30 \pm 7.26	89.45 \pm 1.07	90.44 \pm 1.08	88.01 \pm 1.60	90.46 \pm 1.04	87.91 \pm 1.60	90.68 \pm 1.46	90.19 \pm 1.64
	50%	89.06 \pm 5.17	88.14 \pm 1.16	89.74 \pm 0.74	86.67 \pm 1.37	89.94 \pm 0.59	86.66 \pm 1.34	89.87 \pm 0.68	90.07 \pm 0.96
	30%	87.34 \pm 5.24	84.47 \pm 2.06	88.05 \pm 0.89	83.48 \pm 1.32	89.35 \pm 0.41	83.72 \pm 1.23	88.52 \pm 0.66	89.54 \pm 0.68
	10%	81.75 \pm 13.35	81.38 \pm 3.24	85.82 \pm 1.87	62.52 \pm 6.48	88.56 \pm 0.51	73.57 \pm 1.59	87.28 \pm 0.61	88.26 \pm 1.71
NG9	90%	88.60 \pm 2.88	89.20 \pm 2.97	89.09 \pm 2.86	88.75 \pm 2.71	89.11 \pm 2.92	87.96 \pm 2.64	88.50 \pm 2.58	88.62 \pm 3.14
	70%	87.24 \pm 1.61	88.27 \pm 1.70	88.03 \pm 1.71	87.62 \pm 1.67	88.09 \pm 1.67	87.42 \pm 1.62	87.74 \pm 1.75	88.09 \pm 1.48
	50%	85.68 \pm 3.03	86.11 \pm 1.41	87.21 \pm 1.65	85.94 \pm 1.44	87.37 \pm 1.60	86.06 \pm 1.38	87.22 \pm 1.11	87.56 \pm 1.47
	30%	83.50 \pm 6.71	82.11 \pm 2.54	85.71 \pm 1.54	82.32 \pm 1.46	87.07 \pm 0.88	83.03 \pm 1.30	86.25 \pm 0.81	86.97 \pm 0.97
	10%	80.60 \pm 9.47	80.16 \pm 2.09	84.55 \pm 1.10	68.64 \pm 4.43	86.17 \pm 0.66	72.78 \pm 1.81	85.56 \pm 0.69	86.28 \pm 1.14
Cornell	90%	58.91 \pm 6.25	52.86 \pm 4.17	65.22 \pm 5.56	62.67 \pm 7.04	56.50 \pm 8.58	60.64 \pm 6.65	60.38 \pm 5.78	59.78 \pm 7.22
	70%	59.00 \pm 3.22	51.69 \pm 4.21	63.66 \pm 3.47	60.26 \pm 3.87	57.82 \pm 3.90	57.97 \pm 4.60	59.25 \pm 2.03	61.34 \pm 4.80
	50%	55.40 \pm 5.04	46.74 \pm 4.41	61.70 \pm 2.98	56.88 \pm 3.37	56.05 \pm 2.73	54.53 \pm 5.04	57.09 \pm 3.25	61.25 \pm 4.41
	30%	49.33 \pm 5.91	42.95 \pm 3.30	60.55 \pm 3.78	50.35 \pm 3.25	54.84 \pm 4.37	51.18 \pm 5.44	56.27 \pm 4.12	61.19 \pm 6.19
	10%	44.97 \pm 3.99	41.91 \pm 1.90	58.07 \pm 5.31	42.71 \pm 1.49	57.51 \pm 3.78	47.11 \pm 7.05	58.58 \pm 4.65	58.17 \pm 8.77
Texas	90%	72.11 \pm 4.83	69.11 \pm 4.73	78.29 \pm 4.96	74.04 \pm 6.46	78.99 \pm 4.78	81.34 \pm 6.48	81.72 \pm 4.67	82.29 \pm 4.27
	70%	71.16 \pm 2.09	68.01 \pm 1.90	78.21 \pm 2.46	71.96 \pm 3.94	78.52 \pm 1.83	80.30 \pm 3.62	79.56 \pm 3.24	80.74 \pm 3.07
	50%	68.98 \pm 2.18	65.04 \pm 3.26	77.29 \pm 2.69	69.16 \pm 3.07	76.73 \pm 1.70	78.95 \pm 2.86	77.86 \pm 1.83	79.60 \pm 2.81
	30%	66.70 \pm 2.84	61.27 \pm 4.18	76.11 \pm 2.73	65.56 \pm 2.89	73.79 \pm 2.16	76.50 \pm 2.32	76.63 \pm 1.02	78.16 \pm 1.86
	10%	58.99 \pm 8.20	54.60 \pm 5.22	74.12 \pm 3.88	51.16 \pm 2.46	71.24 \pm 1.78	69.42 \pm 4.04	75.04 \pm 1.53	76.12 \pm 3.36
Washington	90%	68.18 \pm 2.98	63.56 \pm 4.59	70.27 \pm 3.71	66.51 \pm 5.74	61.87 \pm 5.84	61.11 \pm 6.81	58.22 \pm 5.31	64.49 \pm 5.44
	70%	62.28 \pm 6.81	63.65 \pm 1.41	69.29 \pm 2.46	65.37 \pm 3.26	59.78 \pm 2.80	59.74 \pm 3.73	55.45 \pm 4.78	64.05 \pm 3.66
	50%	61.91 \pm 5.08	64.38 \pm 1.15	67.75 \pm 1.97	63.78 \pm 2.36	59.01 \pm 2.21	57.01 \pm 3.97	55.56 \pm 4.18	60.99 \pm 3.70
	30%	60.57 \pm 4.26	64.81 \pm 1.62	65.80 \pm 3.01	59.88 \pm 2.28	57.85 \pm 2.80	53.29 \pm 5.88	55.79 \pm 3.30	60.88 \pm 4.37
	10%	53.52 \pm 12.15	57.73 \pm 10.50	68.04 \pm 1.40	42.18 \pm 4.37	52.11 \pm 5.61	46.30 \pm 8.02	57.30 \pm 3.30	61.57 \pm 5.84
Wisconsin	90%	71.39 \pm 4.90	67.92 \pm 3.40	74.18 \pm 4.71	74.95 \pm 5.03	73.42 \pm 4.10	80.43 \pm 4.82	77.47 \pm 4.95	73.46 \pm 3.62
	70%	70.41 \pm 1.89	67.27 \pm 2.54	75.21 \pm					

TABLE VI

ONE-SIDE t -TEST FOR ALL LABELING RATES. THE WIN/TIE/LOSE FREQUENCY SUMMARIZES HOW MANY TIMES THE BoP CLASSIFIER WAS SIGNIFICANTLY BETTER (WIN), EQUIVALENT (TIE) OR SIGNIFICANTLY WORSE (LOSE) THAN EACH OTHER METHOD.

	RL	RNL	LP	RCT	HF	RWWR	DW1	DW2
$l = 90\%$	6/5/3	9/1/4	8/2/4	3/2/9	9/2/3	5/3/6	12/1/1	6/6/2
$l = 70\%$	14/0/0	11/1/2	10/1/3	3/5/6	11/1/2	6/3/5	10/2/2	9/1/4
$l = 50\%$	13/0/1	13/0/1	11/0/3	10/1/3	12/0/2	9/4/1	12/0/2	9/1/4
$l = 30\%$	13/1/0	13/0/1	13/0/1	10/2/2	14/0/0	10/1/3	14/0/0	12/1/1
$l = 10\%$	14/0/0	14/0/0	14/0/0	9/2/3	14/0/0	7/1/6	14/0/0	11/3/0

TABLE VII

FOR EACH LABELING RATE, THE DIFFERENT CLASSIFIERS ARE RANKED THROUGH A BORDA RATING (SEE THE TEXT FOR DETAILS). THE CLASSIFIERS ARE THEN RANKED ACCORDING TO THE TOTAL RATING OBTAINED ACROSS ALL DATA SETS (THE LARGER THE BETTER). l STANDS FOR LABELING RATE AND THE NUMBERS BETWEEN PARENTHESES ARE THE TOTAL RATINGS.

Ranking	First	Second	Third	Fourth	Fifth	Sixth	Seventh	Last
$l = 90\%$	RCT (86)	BoP (77)	RWWR (76)	DW2 (66)	RL (55)	HF (52)	RNL (49)	DW1 (43)
$l = 70\%$	BoP (87)	RCT (81)	RWWR (76)	DW2 (68)	HF (57)	DW1 (55)	RNL (48)	RL (32)
$l = 50\%$	BoP (94)	RWWR (78)	RCT (75)	DW2 (72)	HF (56)	DW1 (54)	RNL (44)	RL (31)
$l = 30\%$	BoP (103)	RWWR (85)	DW2 (80)	RCT (78)	RNL (45)	DW1 (42)	RL (36)	HF (35)
$l = 10\%$	BoP (100)	RWWR (89)	RCT (85)	DW2 (83)	RNL (48)	RL (42)	DW1 (39)	HF (18)

TABLE VIII

OVERVIEW OF CPU TIME IN SECONDS NEEDED TO CLASSIFY ALL THE UNLABELED NODES. RESULTS ARE AVERAGED ON 20 RUNS. THE CPU USED WAS AN INTEL(R)CORE(TM)I3 AT 2.13 GHZ WITH 3072 Ko OF CACHE SIZE AND 6 GB OF RAM AND THE PROGRAMMING LANGUAGE IS MATLAB.

	RL	RNL	RCT	HF	RWWR	DW1	DW2	BoP
Exp1: 1000 nodes, 2 classes	0.0872	0.4937	0.0751	0.1680	0.3480	0.6826	0.6772	0.4997
Exp2: 2000 nodes, 2 classes	0.4616	3.4961	0.4225	0.9618	2.0441	4.5561	4.5858	3.0574
Exp3: 4000 nodes, 2 classes	2.8274	27.0695	2.5949	7.1481	14.1116	35.7161	36.0207	22.393
Ratio Exp2/Exp1	5.2935	7.0814	5.6258	5.7250	5.8739	6.6746	6.7717	6.1185
Ratio Exp3/Exp2	6.1252	7.7428	6.1418	7.4320	6.9036	7.8392	7.8548	7.3242
Exp2: 2000 nodes, 2 classes	0.4616	3.4961	0.4225	0.9618	2.0441	4.5561	4.5858	3.0574
Exp4: 2000 nodes, 4 classes	0.5011	3.4563	0.4036	1.2064	3.4249	8.5048	8.4003	3.2535
Exp5: 2000 nodes, 8 classes	0.4813	3.8449	0.4482	1.5748	6.0697	16.0031	16.3956	3.5868
Ratio Exp4/Exp2	1.0856	0.9886	0.9553	1.2543	1.6755	1.8667	1.8318	1.0641
Ratio Exp5/Exp4	0.9605	1.1124	1.1105	1.3054	1.7722	1.8817	1.9518	1.1024

majority class with 1768 nodes and the minority class with only 30 nodes (this represents 1.67% for the minority class). Once again, the parameters were tuned by a 10-fold cross-validation and 20 runs were averaged. The classification accuracies for the two classes are reported in Table IX. The best methods to identify the minority class are RL and RNL, followed by RWWR and RCT. In this particular case, the BoP classifier does not perform very well.

VII. CONCLUSION

This paper investigates an application of the bag-of-paths framework viewing the graph as a virtual bag from which paths are drawn according to a Boltzmann sampling distribution.

In particular, it introduces a novel algorithm for graph-based semi-supervised classification through the bag-of-paths group betweenness, or BoP for short (described in Section V). The algorithm sums the a posteriori probabilities of drawing a path visiting a given node of interest according to a biased sampling distribution, and this sum defines our BoP betweenness measure. The Boltzmann sampling distribution depends on a parameter, θ , gradually biasing the distribution towards shorter paths: when θ is large, only little exploration

is performed and only the shortest paths are considered, while when θ is small (close to 0^+), longer paths are considered and are sampled according to the product of the transition probabilities p_{ij}^{ref} along the path (a natural random walk).

Experiments on real-world data sets show that the BoP method outperforms the other considered approaches when only a few labeled nodes are available. When more nodes are labeled, the BoP method is still competitive. Its computation time is also substantially lower in most of the cases.

Our future work will include several extensions of the proposed approach. Another interesting issue is how to combine the information provided by the graph and the features of the nodes in a clever, preferably optimal, way. The interest of including node features should be assessed experimentally. A typical case study could be the labeling of protein-protein interaction networks. The node features could involve gene expression measurements for the corresponding proteins.

Yet another application of the bag-of-paths framework could be the definition of a robustness measure or criticality measure of the nodes. The idea would be to compute the change in reachability between nodes when deleting one node within the BoP framework. Nodes having a large impact on reachability would be then considered as highly critical.

TABLE IX

CLASSIFICATION ACCURACIES IN PERCENTS, AVERAGED OVER 20 RUNS, OBTAINED ON PARTIALLY LABELED ARTIFICIAL GRAPHS. RESULTS ARE REPORTED FOR THE EIGHT METHODS (RL, NRL, RCT, HF, RWWR, DW1, DW2, BoP) AND FOR A 50% LABELING RATE. RARE IS THE CASE WHERE ONLY TWO LABELED NODES PER CLASS ARE KNOWN AND IMBALANCED IS THE CASE WHERE ONE OF THE CLASSES IS MUCH MORE REPRESENTED THAN THE OTHER.

	RL	NRL	RCT	HF	RWWR	DW1	DW2	BoP
Rare	51.9 \pm 1.4	52.6 \pm 2.6	84.0 \pm 2.5	50.0 \pm 0.05	84.0 \pm 2.5	51.9 \pm 0.5	67.3 \pm 2.8	83.0 \pm 3.1
Imbalanced: Major	98.2 \pm 0.04	98.2 \pm 0.04	98.1 \pm 0.08	99.9 \pm 0.04	95.5 \pm 0.03	88.7 \pm 7.0	93.8 \pm 2.5	96.8 \pm 0.4
Imbalanced: Minor	42.4 \pm 4.2	43.8 \pm 3.3	32.2 \pm 2.0	1.9 \pm 2.4	11.6 \pm 10.5	38.5 \pm 0.0	17.2 \pm 2.7	11.9 \pm 2.5

Finally, the biggest drawback of the BoP classifier is that it is not applicable as-is for large graphs. It would be interesting to investigate if it is possible to modify the classifier to be computationally more tractable on large graphs. A starting clue would be to use the same trick as in [8].

REFERENCES

- [1] K. Francoise, I. Kivimaki, A. Mantrach, F. Rossi, and M. Saerens, "A bag-of-paths framework for network data analysis," *Submitted for publication and available on ArXiv at <http://arxiv.org/abs/1302.6766>*.
- [2] J. Callut and P. Dupont, "Learning hidden markov models from first passage times," *Proceedings of the European Machine Learning conference (ECML). Lecture notes in Artificial Intelligence, Springer, 2007*.
- [3] J.-Y. Pan, H.-J. Yang, C. Faloutsos, and P. Duygulu, "Automatic multimedia cross-modal correlation discovery," *Proceedings of the 10th ACM SIGKDD international conference on Knowledge Discovery and Data Mining (KDD 2004)*, pp. 653–658, 2004.
- [4] D. Zhou, O. Bousquet, T. Lal, J. Weston, and B. Scholkopf, "Learning with local and global consistency," in *Proceedings of the Neural Information Processing Systems Conference (NIPS 2003)*, 2003, pp. 237–244.
- [5] X. Zhu, Z. Ghahramani, and J. Lafferty, "Semi-supervised learning using gaussian fields and harmonic functions," in *Proceedings of the Twentieth International Conference on Machine Learning (ICML 2003)*, 2003, pp. 912–919.
- [6] X. Zhu and A. Goldberg, *Introduction to semi-supervised learning*. Morgan & Claypool Publishers, 2009.
- [7] A. Mantrach, L. Yen, J. Callut, K. Francoise, M. Shimbo, and M. Saerens, "The sum-over-paths covariance kernel: a novel covariance between nodes of a directed graph," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 6, pp. 1112–1126, 2010.
- [8] A. Mantrach, N. van Zeebroeck, P. Francq, M. Shimbo, H. Bersini, and M. Saerens, "Semi-supervised classification and betweenness computation on large, sparse, directed graphs," *Pattern Recognition*, vol. 44, no. 6, pp. 1212 – 1224, 2011.
- [9] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon, "Network motifs: Simple building blocks of complex networks," *Science*, vol. 298, p. p.824, 2002.
- [10] A. Arenas, A. Fernandez, S. Fortunato, and S. Gomez, "Motif-based communities in complex networks," *Journal of Physics A: Mathematical and Theoretical*, vol. 41, p. 224001, 2008.
- [11] S. Abney, *Semisupervised learning for computational linguistics*. Chapman and Hall/CRC, 2008.
- [12] O. Chapelle, B. Scholkopf, and A. Zien (editors), *Semi-supervised learning*. MIT Press, 2006.
- [13] D. Zhou and B. Scholkopf, "Learning from labeled and unlabeled data using random walks," in *Proceedings of the 26th DAGM Symposium, 2004*, pp. 237–244.
- [14] M. Szummer and T. Jaakkola, "Partially labeled classification with markov random walks," in *Advances in Neural Information Processing Systems*, T. Dietterich, S. Becker, and Z. Ghahramani, Eds., vol. 14. Vancouver, Canada: MIT Press, 2001.
- [15] O. Chapelle, J. Weston, and B. Scholkopf, "Cluster kernels for semi-supervised learning," in: *NIPS, 2002*, pp. 585–592, 2002.
- [16] A. Kapoor, Y. Qi, H. Ahn, and R. Picard, "Hyperparameter and kernel learning for graph based semi-supervised classification," in: *NIPS, 2005*, pp. 627–634, 2002.
- [17] D. Zhou, J. Huang, and B. Scholkopf, "Learning from labeled and unlabeled data on a directed graph," in *Proceedings of the 22nd International Conference on Machine Learning, 2005*, pp. 1041–1048.
- [18] M. Belkin, I. Matveeva, and P. Niyogi, "Regularization and semi-supervised learning on large graphs," in *Proceedings of the International Conference on Learning Theory (COLT 2004)*, 2004, pp. 624–638.
- [19] J. Wang, F. Wang, C. Zhang, H. Shen, and L. Quan, "Linear neighborhood propagation and its applications," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 9, pp. 1600–1615, 2009.
- [20] T. Joachims, "Transductive learning via spectral graph partitioning," in *Proceedings of the 20th International Conference on Machine Learning (ICDM 2003)*, Washington DC, 2003, p. 290–297.
- [21] T. Kato, H. Kashima, and M. Sugiyama, "Robust label propagation on multiple networks," *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 35–44, 2009.
- [22] F. Fouss, K. Francoise, L. Yen, A. Pirotte, and M. Saerens, "An experimental investigation of kernels on a graph on collaborative recommendation and semisupervised classification," *Neural Networks*, vol. 31, pp. 53–72, 2012.
- [23] D. Zhou and B. Scholkopf, "Discrete regularization," in *Semi-supervised learning, O. Chapelle, B. Scholkopf and A. Zien (editors)*. MIT Press, 2006, pp. 237–249.
- [24] S. A. Macskassy and F. Provost, "Classification in networked data: a toolkit and a univariate case study," *Journal of Machine Learning Research*, vol. 8, pp. 935–983, 2007.
- [25] L. Getoor and B. Taskar (editors), *Introduction to statistical relational learning, 2007*.
- [26] H. Tong, C. Faloutsos, and J.-Y. Pan, "Random walk with restart: fast solutions and applications," *Knowledge and Information Systems*, vol. 14, no. 3, pp. 327–346, 2008.
- [27] D. Liben-Nowell and J. Kleinberg, "The link-prediction problem for social networks," *Journal of the American Society for Information Science and Technology*, vol. 58, no. 7, pp. 1019–1031, 2007.
- [28] M. Brand and K. Huang, "A unifying theorem for spectral embedding and clustering," in *Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics*, Key West, FL, January 2003.
- [29] P. Sarkar and A. Moore, "A tractable approach to finding closest truncated-commute-time neighbors in large graphs," *Proceedings of the 23rd Conference on Uncertainty in Artificial Intelligence (UAI)*, 2007.
- [30] J. Callut, K. Francoise, M. Saerens, and P. Dupont, "Semi-supervised classification from discriminative random walks," *Proceedings of the European Machine Learning conference (ECML 2008). Lecture notes in Artificial Intelligence, Springer, 2008*, pp. 162–177, 2008.
- [31] M. Herbster, M. Pontil, and S. Rojas-Galeano, "Fast prediction on tree," *Proceedings of the 22th Neural Information Processing Conference NIPS 2008*, pp. 657–664, 2008.
- [32] L. Tang and H. Liu, "Relational learning via latent social dimensions," in *Proceedings of the ACM conference on Knowledge Discovery and Data Mining (KDD 2009)*, 2009, pp. 817–826.
- [33] —, "Scalable learning of collective behavior based on sparse social dimensions," in *Proceedings of the ACM conference on Information and Knowledge Management (CIKM 2009)*, 2009, pp. 1107–1116.
- [34] —, "Toward predicting collective behavior via social dimension extraction," *IEEE Intelligent Systems*, vol. 25, no. 4, pp. 19–25, 2010.
- [35] L. Gettor and B. Taskar, Eds., *Introduction to Statistical Relational Learning*. The MIT Press, 2007.
- [36] L. Freeman, "A set of measures of centrality based on betweenness," *Sociometry*, vol. 40, no. 1, pp. 35–41, 1977.
- [37] M. Newman, "A measure of betweenness centrality based on random walks," *Social Networks*, vol. 27, no. 1, pp. 39–54, 2005.
- [38] E. D. Kolaczyk, *Statistical analysis of network data: methods and models*. Springer, 2009.
- [39] K. Lang, "Newsweeder: Learning to filter netnews," in *Proceedings of the Twelfth International Conference on Machine Learning, 1995*, pp. 331–339.

- [40] L. Yen, F. Fouss, C. Decaestecker, P. Francq, and M. Saerens, "Graph nodes clustering with sigmoid commute-time kernel : A comparative study," *Data & Knowledge Engineering*, no. 68, pp. 338–361, 2009.