# AdaDIF: Adaptive Diffusions for Efficient Semi-supervised Learning over Graphs

Dimitris Berberidis<sup>1</sup>, Athanasios N. Nikolakopoulos<sup>2</sup> and Georgios B. Giannakis<sup>1,2</sup>

<sup>1</sup>Department of Electrical & Computer Engineering

<sup>2</sup>Digital Technology Center, University of Minnesota, Minneapolis, MN, USA

Emails: {bermp,anikolak,georgios}@umn.edu

Abstract-Diffusion-based classifiers such as those relying on the Personalized PageRank and the Heat kernel, enjoy remarkable classification accuracy at modest computational requirements. Their performance however is affected by the extent to which the chosen diffusion captures a typically unknown label propagation mechanism, that can be specific to the underlying graph, and potentially different for each class. The present work introduces a disciplined, data-efficient approach to learning class-specific diffusion functions adapted to the underlying network topology. The novel learning approach leverages the notion of "landing probabilities" of class-specific random walks, which can be computed efficiently, thereby ensuring scalability to large graphs. This is supported by rigorous analysis of the properties of the model as well as the proposed algorithms. Classification tests on real networks demonstrate that adapting the diffusion function to the given graph and observed labels, significantly improves the performance over fixed diffusions; reaching—and many times surpassing—the classification accuracy of computationally heavier state-of-theart competing methods, that rely on node embeddings and deep neural networks.

Keywords-Random Walks; Networks; Markov Chains; Label Propagation; Dictionary

## I. INTRODUCTION

The task of classifying nodes of a graph arises frequently in several applications on real-world networks, such as the ones derived from social interactions and biological dependencies. Graph-based semi-supervised learning (SSL) methods tackle this task building on the premise that the true labels are distributed "smoothly" with respect to the underlying network, which then motivates leveraging the network structure to increase the classification accuracy [11]. Graph-based SSL has been pursued by various intertwined methods, including iterative label propagation [6], [38], [22], kernels on graphs [27], manifold regularization [5], graph partitioning [35], [17], transductive learning [34], competitive infection models [32], and bootstrapped label propagation [10]. Recently, approaches based on node-embeddings [30], [16], [37], as well as deep-learning architectures [18], [2] have gained popularity, and were reported to have state-ofthe-art performance.

Many of the aforementioned methods are challenged by computational complexity and scalability issues that limit their applicability to large-scale networks. Random-walk-based diffusions present an efficient and effective alternative. Methods of this family diffuse probabilistically the known labels through the graph, thereby ranking nodes according to weighted sums of variable-length landing probabilities. Celebrated representatives include those based on the Personalized PageRank and the Heat Kernel that were found to perform remarkably well in certain application domains [19], and have been nicely linked to particular network models [20], [3], [21]. However, the effectiveness of diffusion-based classifiers can vary considerably depending on whether the chosen diffusion conforms with the latent label propagation mechanism that might be, (i) particular to the target application or underlying network topology; and, (ii) different for each class.

The present contribution alleviates these shortcomings and markedly improves the performance of random-walk-based classifiers by *adapting the diffusion functions* to both the network and the observed labels. The resulting novel classifier relies on the notion of landing probabilities of *short-length random walks* rooted at the observed nodes of each class. The small number of these landing probabilities can be extracted efficiently with a small number of sparse matrix-vector products, thus ensuring applicability to large-scale networks<sup>1</sup>. Theoretical analysis establishes that short random walks are in most cases sufficient for reliable classification. We test our methods in terms of multiclass and multilabel classification accuracy, and confirm that it can achieve results competitive to state-of-the-art methods, while also being considerably faster.

# II. PROBLEM STATEMENT AND MODELING

Consider a graph  $\mathcal{G} := \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V}$  is the set of N nodes, and  $\mathcal{E}$  the set of edges. Connectivity is captured by the weight matrix  $\mathbf{W}$  having entries  $W_{ij} > 0$  if  $(i,j) \in \mathcal{E}$ . Associated with each node  $i \in \mathcal{V}$  there is a discrete label  $y_i \in \mathcal{Y}$ . In SSL classification over graphs, a subset  $\mathcal{L} \subset \mathcal{V}$  of nodes has available labels  $\mathbf{y}_{\mathcal{L}}$ , and the goal is to infer the labels of the unlabeled set  $\mathcal{U} := \mathcal{V} \setminus \mathcal{L}$ . Given a measure of influence, a node most influenced by labeled nodes of a certain class is deemed to also belong to the same class. Thus, label-propagation on graphs boils down to quantifying the influence of  $\mathcal{L}$  on  $\mathcal{U}$ , see, e.g. [11], [22], [36]. An intuitive

<sup>&</sup>lt;sup>1</sup>This work was supported by NSF 1711471, 1500713, and 1442686.

<sup>&</sup>lt;sup>1</sup>Scalable implementation available at: github.com/DimBer/SSL\_lib

yet simple measure of node-to-node influence relies on the notion of random walks on graphs.

A simple random walk on a graph is a discrete-time Markov chain with state space the set of nodes, and transition probabilities

$$[\mathbf{H}]_{ij} := \Pr\{X_k = i | X_{k-1} = j\} = W_{ij}/d_j = [\mathbf{W}\mathbf{D}^{-1}]_{ij}$$

where  $X_k \in \mathcal{V}$  denotes the position of the random walker (state) at the  $k^{th}$  step;  $d_j := \sum_{k \in \mathcal{N}_j} W_{kj}$  is the degree of node j; and,  $\mathcal{N}_j$  its neighborhood. Since we consider undirected graphs the steady-state distribution of  $\{X_k\}$  always exists if it is connected, and non-bipartite. It is given by the dominant right eigenvector of the column-stochastic transition probability matrix  $\mathbf{H} := \mathbf{W}\mathbf{D}^{-1}$ , where  $\mathbf{D} := \mathrm{diag}\,(d_1,d_2,\ldots,d_N)$  [24]. The steady-state distribution  $\pi$  can be shown to have entries

$$\pi_i := \lim_{k \to \infty} \sum_{i \in \mathcal{V}} \Pr\{X_k = i | X_0 = j\} \Pr\{X_0 = j\} = \frac{d_i}{2|\mathcal{E}|}$$

that are clearly not dependent on the initial "seeding" distribution  $\Pr\{X_0\}$ ; and  $\pi$  is thus unsuitable for measuring influence among nodes. Instead, for graph-based SSL, we will utilize the k-step landing probability per node i given by

$$p_i^{(k)} := \sum_{j \in \mathcal{V}} \Pr\{X_k = i | X_0 = j\} \Pr\{X_0 = j\}$$
 (1)

that in vector form  $\mathbf{p}^{(k)} := [p_1^{(k)} \dots p_N^{(k)}]^\mathsf{T}$  satisfies  $\mathbf{p}^{(k)} = \mathbf{H}^k \mathbf{p}^{(0)}$ , where  $p_i^{(0)} := \Pr\{X_0 = i\}$ . In words,  $p_i^{(k)}$  is the probability that a random walker with initial distribution  $\mathbf{p}^{(0)}$  is located at node i after k steps. Therefore,  $p_i^{(k)}$  is a valid measure of the influence that  $\mathbf{p}^{(0)}$  has on any node in  $\mathcal{V}$ .

The landing probabilities per class  $c \in \mathcal{Y}$  are (cf. (1))

$$\mathbf{p}_c^{(k)} = \mathbf{H}^k \mathbf{v}_c \tag{2}$$

where for  $\mathcal{L}_c := \{i \in \mathcal{L} : y_i = c\}$ , we select as  $\mathbf{v}_c$  the normalized class-indicator vector with i-th entry

$$[\mathbf{v}_c]_i = \begin{cases} 1/|\mathcal{L}_c|, & i \in \mathcal{L}_c \\ 0, & \text{else} \end{cases}$$
 (3)

acts as initial distribution. Using (2), we model diffusions per class c over the graph driven by  $\{\mathbf{p}_c^{(k)}\}_{k=1}^K$  as

$$\mathbf{f}_c(\boldsymbol{\theta}) = \sum_{k=1}^K \theta_k \mathbf{p}_c^{(k)} = \mathbf{P}_c^{(K)} \boldsymbol{\theta}$$
 (4)

where  $\mathbf{P}_c^{(K)} := \begin{bmatrix} \mathbf{p}_c^{(1)} & \cdots & \mathbf{p}_c^{(K)} \end{bmatrix}$ , and  $\theta_k$  denotes the importance assigned to the  $k^{th}$  hop neighborhood. By constraining  $\boldsymbol{\theta} \in \mathcal{S}^K$ , where  $\mathcal{S}^K := \{\mathbf{x} \in \mathbb{R}^K : \mathbf{x} \geq \mathbf{0}, \mathbf{1}^\mathsf{T}\mathbf{x} = 1\}$  is the K-dimensional probability simplex,  $\mathbf{f}_c(\boldsymbol{\theta})$  becomes a valid nodal probability mass function (pmf) for class c.

Given  $\theta$  and upon obtaining  $\{\mathbf{f}_c(\theta)\}_{c\in\mathcal{Y}}$ , our diffusion-based classifiers will predict labels over  $\mathcal{U}$  as

$$\hat{y}_i(\boldsymbol{\theta}) := \arg \max_{c \in \mathcal{Y}} \left[ \mathbf{f}_c(\boldsymbol{\theta}) \right]_i \tag{5}$$

where  $[\mathbf{f}_c(\boldsymbol{\theta})]_i$  is the  $i^{th}$  entry of  $\mathbf{f}_c(\boldsymbol{\theta})$ .

Next, we outline two notable members of the family of diffusion-based classifiers that can be viewed as special cases of (4).

#### A. Personalized-PageRank and Heat-Kernel Classifiers

Inspired by its celebrated network centrality metric [9], the Personalized PageRank (PPR) algorithm has well-documented merits for label propagation; see, e.g. [25]. PPR is a special case of (4) corresponding to  $\theta_{\rm PPR} = (1-\alpha)\left[\alpha,\alpha^2,\ldots,\alpha^K\right]^{\sf T}$ , where  $0<\alpha<1$ , and  $1-\alpha$  can be interpreted as the "restart" probability of random walks with restarts.

The PPR-based classifier relies on (cf. (4))

$$\mathbf{f}_c(\boldsymbol{\theta}_{PPR}) = (1 - \alpha) \sum_{k=0}^{K} \alpha^k \mathbf{p}_c^{(k)}$$
 (6)

satisfying asymptotically in the number of random walk steps

$$\lim_{K \to \infty} \mathbf{f}_c(\boldsymbol{\theta}_{PPR}) = (1 - \alpha)(\mathbf{I} - \alpha \mathbf{H})^{-1} \mathbf{v}_c$$

which implies that  $\mathbf{f}_c(\boldsymbol{\theta}_{\mathrm{PPR}})$  approximates the solution of a linear system. Indeed, as shown in [3], PPR amounts to solving a weighted regularized least-squares problem over  $\mathcal{V}$ ; see also [20] for a PPR interpretation as an approximate geometric discriminant function defined in the space of landing probabilities.

The heat kernel (HK) is another popular diffusion that has recently been employed for SSL [27] and community detection on graphs [19]. HK is also a special case of (4) with  $\boldsymbol{\theta}_{\rm HK} = e^{-t} \left[t, \frac{t^2}{2}, \dots, \frac{t^K}{K!}\right]^{\sf T}$ , yielding class distributions (cf. (4))

$$\mathbf{f}_c(\boldsymbol{\theta}_{\mathrm{HK}}) = e^{-t} \sum_{k=0}^{K} \frac{t^k}{k!} \mathbf{p}_c^{(k)}.$$
 (7)

Furthermore, it can be readily shown that

$$\lim_{K \to \infty} \mathbf{f}_c(\boldsymbol{\theta}_{\mathrm{HK}}) = e^{-t(\mathbf{I} - \mathbf{H})} \mathbf{v}_c$$

allowing HK to be interpreted as an approximation of a heat diffusion process, where heat is flowing from  $\mathcal{L}_c$  to the rest of the graph; and  $\mathbf{f}_c(\boldsymbol{\theta}_{\rm HK})$  is a snapshot of the temperature after time t has elapsed. HK provably yields low conductance communities, while also converging faster to its asymptotic closed-form expression than PPR [13].

#### III. ADAPTIVE DIFFUSIONS

Besides the unifying view of (4), the main contribution here is on efficiently designing  $\mathbf{f}_c(\theta_c)$  in (4), by learning the corresponding  $\theta_c$  per class. Thus, unlike PPR and HK, the methods introduced here can afford class-specific label propagation that is *adaptive* to the graph structure, and also *adaptive* to the labeled nodes; see Fig. 1 for a high-level illustration of the proposed adaptive diffusion framework.

Consider for generality a goodness-of-fit loss  $\ell(\cdot)$ , and a regularizer  $R(\cdot)$  promoting e.g., smoothness over the graph. Using these, the sought class distribution will be

$$\hat{\mathbf{f}}_c = \arg\min_{\mathbf{f} \in \mathbb{R}^N} \ell(\mathbf{y}_{\mathcal{L}_c}, \mathbf{f}) + \lambda R(\mathbf{f})$$
 (8)

where  $\lambda$  tunes the degree of regularization, and

$$[\mathbf{y}_{\mathcal{L}_c}]_i = \begin{cases} 1, & i \in \mathcal{L}_c \\ 0, & \text{else} \end{cases}$$

is the indicator vector of the nodes belonging to class c. Using our diffusion model in (4), the N-dimensional optimization problem (8) reduces to solving for the K-dimensional vector ( $K \ll N$ )

$$\hat{\boldsymbol{\theta}}_c = \arg\min_{\boldsymbol{\theta} \in \mathcal{S}^K} \ell(\mathbf{y}_{\mathcal{L}_c}, \mathbf{f}_c(\boldsymbol{\theta})) + \lambda R(\mathbf{f}_c(\boldsymbol{\theta})). \tag{9}$$

Although many choices of  $\ell(\cdot)$  may be of interest, we will focus for simplicity on the quadratic loss, namely

$$\ell(\mathbf{y}_{\mathcal{L}_c}, \mathbf{f}) := \sum_{i \in \mathcal{L}} \frac{1}{d_i} ([\bar{\mathbf{y}}_{\mathcal{L}_c}]_i - f_i)^2$$
$$= (\bar{\mathbf{y}}_{\mathcal{L}_c} - \mathbf{f})^\mathsf{T} \mathbf{D}_c^\dagger (\bar{\mathbf{y}}_{\mathcal{L}_c} - \mathbf{f})$$
(10)

where  $\bar{\mathbf{y}}_{\mathcal{L}_c} := (1/|\mathcal{L}|) \, \mathbf{y}_{\mathcal{L}_c}$  is the class indicator vector after normalization to avoid overfitting and numerical instabilities, and  $\mathbf{D}_{\mathcal{L}}^{\dagger} = \operatorname{diag}(\mathbf{d}_{\mathcal{L}}^{(-1)})$  with entries

$$[\mathbf{d}_{\mathcal{L}}^{(-1)}]_i = \begin{cases} 1/d_i, & i \in \mathcal{L} \\ 0, & \text{else} \end{cases}.$$

For a smoothness-promoting regularization, we will employ the following (normalized) Laplacian-based metric

$$R(\mathbf{f}) = \frac{1}{2} \sum_{i \in \mathcal{V}} \sum_{i \in \mathcal{N}_i} \left( \frac{f_i}{d_i} - \frac{f_j}{d_j} \right)^2 = \mathbf{f}^\mathsf{T} \mathbf{D}^{-1} \mathbf{L} \mathbf{D}^{-1} \mathbf{f}. \quad (11)$$

Intuitively speaking, (10) favors vectors  $\mathbf{f}$  having non-zero  $(|1/|\mathcal{L}|)$  values on nodes that are known to belong to class c, and zero values on nodes that are known to belong to other classes  $(\mathcal{L} \setminus \mathcal{L}_c)$ , while (11) promotes similarity of the entries of  $\mathbf{f}$  that correspond to neighboring nodes. In (10) and (11), each entry  $f_i$  is normalized by  $d_i^{-\frac{1}{2}}$  and  $d_i^{-1}$  respectively. This normalization counterbalances the tendency of random walks to concentrate on high-degree nodes, thus placing equal importance to all nodes.

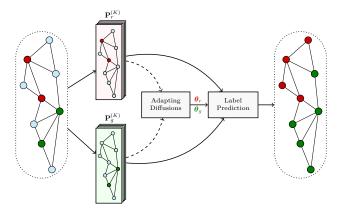


Figure 1. High-level illustration of adaptive diffusions. The nodes belong to two classes (**red** and **green**). The per-class diffusions are learned by exploiting the landing probability spaces produced by random walks rooted at the sample nodes (second layer: **up** for red; **down** for green).

Substituting (10) and (11) into (9), and recalling from (4) that  $\mathbf{f}_c(\boldsymbol{\theta}) = \mathbf{P}_c^{(K)} \boldsymbol{\theta}$ , yields the convex quadratic program

$$\hat{\boldsymbol{\theta}}_c = \arg\min_{\boldsymbol{\theta} \in S^K} \boldsymbol{\theta}^\mathsf{T} \mathbf{A}_c \boldsymbol{\theta} + \boldsymbol{\theta}^\mathsf{T} \mathbf{b}_c \tag{12}$$

with  $\mathbf{b}_c$  and  $\mathbf{A}_c$  given by

$$\mathbf{b}_c = -\frac{2}{|\mathcal{L}|} (\mathbf{P}_c^{(K)})^\mathsf{T} \mathbf{D}_{\mathcal{L}}^\dagger \mathbf{y}_{\mathcal{L}_c}$$
 (13)

$$\mathbf{A}_{c} = (\mathbf{P}_{c}^{(K)})^{\mathsf{T}} \mathbf{D}_{\mathcal{L}}^{\dagger} \mathbf{P}_{c}^{(K)} + \lambda (\mathbf{P}_{c}^{(K)})^{\mathsf{T}} \mathbf{D}^{-1} \mathbf{L} \mathbf{D}^{-1} \mathbf{P}_{c}^{(K)}$$
(14)

$$= (\mathbf{P}_{c}^{(K)})^{\mathsf{T}} \left[ \left( \mathbf{D}_{\mathcal{L}}^{\dagger} + \lambda \mathbf{D}^{-1} \right) \mathbf{P}_{c}^{(K)} - \lambda \mathbf{D}^{-1} \mathbf{H} \mathbf{P}_{c}^{(K)} \right]$$
$$= (\mathbf{P}_{c}^{(K)})^{\mathsf{T}} \left( \mathbf{D}_{\mathcal{L}}^{\dagger} \mathbf{P}_{c}^{(K)} + \lambda \mathbf{D}^{-1} \tilde{\mathbf{P}}_{c}^{(K)} \right)$$
(15)

where

$$\begin{aligned} \mathbf{H}\mathbf{P}_{c}^{(K)} &= \begin{bmatrix} \mathbf{H}\mathbf{p}_{c}^{(1)} & \mathbf{H}\mathbf{p}_{c}^{(2)} & \cdots & \mathbf{H}\mathbf{p}_{c}^{(K)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{p}_{c}^{(2)} & \mathbf{p}_{c}^{(3)} & \cdots & \mathbf{p}_{c}^{(K+1)} \end{bmatrix} \end{aligned}$$

is a "shifted" version of  $\mathbf{P}_c^{(K)}$ , where each  $\mathbf{p}_c^{(k)}$  is advanced by one step, and

$$ilde{\mathbf{P}}_c^{(K)} := egin{bmatrix} ilde{\mathbf{p}}_c^{(1)} & ilde{\mathbf{p}}_c^{(2)} & \cdots & ilde{\mathbf{p}}_c^{(K)} \end{bmatrix}$$

with  $\tilde{\mathbf{p}}_c^{(i)} := \mathbf{p}_c^{(i)} - \mathbf{p}_c^{(i+1)}$  containing the "differential" landing probabilities. The complexity of "naively" finding the  $K \times K$  matrix  $\mathbf{A}_c$  (and thus also  $\mathbf{b}_c$ ) is  $\mathcal{O}(K^2N)$  for computing the first summand, and  $\mathcal{O}(|\mathcal{E}|K)$  for the second summand in (14), after leveraging the sparsity of  $\mathbf{L}$ , which means  $|\mathcal{E}| \ll N^2$ . But since columns of  $\tilde{\mathbf{P}}_c^{(K)}$  are obtained as differences of consecutive columns of  $\mathbf{P}_c^{(K)}$ , this load of  $\mathcal{O}(|\mathcal{E}|K)$  is saved. In a nutshell, the adaptive-diffusion (AdaDIF) solver in (12)-(15) incurs complexity  $\mathcal{O}(K^2N)$ .

### A. Limiting behavior and computational complexity

In this section, we offer further insights on the model (4), along with complexity analysis of the parametric solution in (12). To start, the next proposition establishes the limiting

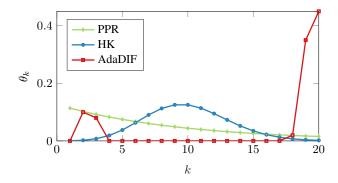


Figure 2. Illustration of K=20 landing probability coefficients for PPR with  $\alpha=0.9$ , HK with t=10, and AdaDIF.

behavior of AdaDIF as the regularization parameter grows; for the proof see [7].

**Proposition 1.** If the second largest eigenvalue of **H** has multiplicity 1, then for K sufficiently large but finite, the solution to (12) as  $\lambda \to \infty$  satisfies

$$\hat{\boldsymbol{\theta}}_c = \mathbf{e}_K, \quad \forall \ \mathcal{L}_c \subseteq \mathcal{V}.$$
 (16)

Our experience with solving (12) reveal that the sufficiently large K required for (16) to hold, can be as small as  $10^2$ . As  $\lambda \to \infty$ , the effect of the loss in (10) vanishes. According to Proposition 1, this causes AdaDIF to boost smoothness by concentrating the simplex weights (entries of  $\hat{\theta}_c$ ) on landing probabilities of the late steps (k close to K). If on the other extreme, smoothness-over-the-graph is not promoted (cf.  $\lambda = 0$ ), the sole objective of AdaDIF is to construct diffusions that best fit the available labeled data. Since short-length random walks from a given node typically lead to nodes of the same class, while longer walks to other classes, AdaDIF with  $\lambda = 0$  tends to leverage only a few landing probabilities of early steps (k close to 1). For moderate values of  $\lambda$ , AdaDIF effectively adapts per-class diffusions by balancing the emphasis on initial versus final landing probabilities. Fig. 2 depicts an example of how AdaDIF places weights  $\{\theta_k\}_{k=1}^K$ on landing probabilities after a maximum of K = 20 steps, generated from few samples belonging to one of 7 classes of the Cora citation network. Note that the learned coefficients may follow radically different patterns than those dictated by standard non-adaptive diffusions such as PPR or HK. It is worth noting that the simplex constraint induces sparsity of the solution in (12), thus 'pushing'  $\{\theta_k\}$  entries to zero.

The computational core of the proposed method is to build the landing probability matrix  $\mathbf{P}_c^{(K)}$ , whose columns are computed fast using power iterations leveraging the sparsity of  $\mathbf{H}$  (cf. (2)). This endows AdaDIF with high computational efficiency, especially for small K. Specifically, since for solving (12) AdaDIF incurs complexity  $\mathcal{O}(K^2N)$  per class, if  $K < |\mathcal{E}|/N$ , this becomes  $\mathcal{O}(|\mathcal{E}|K)$ ; and for  $|\mathcal{Y}|$  classes, the overall complexity of AdaDIF is  $\mathcal{O}(|\mathcal{Y}||\mathcal{E}|K)$ , which is in the same order as that of non-adaptive diffusions such as

PPR and HK. For larger K however, an additional  $\mathcal{O}(K^2N)$  is required per class, mainly to obtain  $\mathbf{A}_c$  in (15).

Overall, if  $\mathcal{O}(KN)$  memory requirements are met, the runtime of AdaDIF scales *linearly* with  $|\mathcal{E}|$ , provided that K remains small. Thankfully, small values of K are usually sufficient to achieve high learning performance. As will be shown in the next section, this observation is in par with the analytical properties of diffusion based classifiers, where it turns out that K large does not improve classification accuracy.

# B. On the choice of K

Here we elaborate on how the selection of K influences the classification task at hand. As expected, the effect of K is intimately linked to the topology of the underlying graph, the labeled nodes, and their properties. For simplicity, we will focus on binary classification with the two classes denoted by "+" and "-." Central to our subsequent analysis is a concrete measure of the effect an extra landing probability vector  $\mathbf{p}_c^{(k)}$  can have on the outcome of a diffusion-based classifier. Intuitively, this effect is diminishing as the number of steps K grows, as both random walks eventually converge to the same stationary distribution. Motivated by this, we introduce next the  $\gamma$ -distinguishability threshold.

**Definition 1** ( $\gamma$ -distinguishability threshold). Let  $\mathbf{p}_+$  and  $\mathbf{p}_-$  denote respectively the seed vectors for nodes of class "+" and "-," initializing the landing probability vectors in matrices  $\mathbf{X}_c := \mathbf{p}_c^{(K)}$ , and  $\check{\mathbf{X}}_c := \left[\mathbf{p}_c^{(1)} \cdots \mathbf{p}_c^{(K-1)} \mathbf{p}_c^{(K+1)}\right]$ , where  $c \in \{+, -\}$ . With  $\mathbf{y} := \mathbf{X}_+ \boldsymbol{\theta} - \mathbf{X}_- \boldsymbol{\theta}$  and  $\check{\mathbf{y}} := \check{\mathbf{X}}_+ \boldsymbol{\theta} - \check{\mathbf{X}}_- \boldsymbol{\theta}$ , the  $\gamma$ -distinguishability threshold of the diffusion-based classifier is the smallest integer  $K_\gamma$  satisfying  $\|\mathbf{y} - \check{\mathbf{y}}\| \leq \gamma$ .

The following theorem establishes an upper bound on  $K_{\gamma}$  expressed in terms of fundamental quantities of the graph, as well as basic properties of the labeled nodes per class; for the proof see [7].

**Theorem 1.** For any diffusion-based classifier with coefficients  $\theta$  constrained to a probability simplex of appropriate dimensions, the  $\gamma$ -distinguishability threshold is upperbounded as

$$K_{\gamma} \leq \frac{1}{\mu'} \log \left[ \frac{2\sqrt{d_{\max}}}{\gamma} \left( \sqrt{\frac{1}{d_{\min}|\mathcal{L}_{-}|}} + \sqrt{\frac{1}{d_{\min}|\mathcal{L}_{+}|}} \right) \right]$$

where  $d_{\min +} := \min_{i \in \mathcal{L}_+} d_i$ ,  $d_{\min -} := \min_{j \in \mathcal{L}_-} d_j$ ,  $d_{\max} := \max_{i \in \mathcal{V}} d_i$  and  $\mu' := \min\{\mu_2, 2 - \mu_N\}$  where  $\{\mu_n\}_{n=1}^N$  denote the eigenvalues of the normalized graph Laplacian in ascending order.

The  $\gamma$ -distinguishability threshold can guide the choice of the dimension K of the landing probability space. Indeed, using class-specific landing probability steps  $K \geq K_{\gamma}$ , does not help distinguishing between the corresponding classes, in the sense that the classifier output is not perturbed by more than  $\gamma$ . Intuitively, the information contained in the landing

probabilities  $K_{\gamma} + 1, K_{\gamma} + 2,...$  is essentially the same for both classes and thus, using them as features unnecessarily increases the overall complexity of the classifier, and also "opens the door" to curse of dimensionality related concerns.

Theorem 1 makes no assumptions on the diffusion coefficients, so long they belong to a probability simplex. Of course, specifying the diffusion function can specialize and further tighten the corresponding  $\gamma$ -distinguishability threshold. Conveniently, our experiments suggest that  $K \in [10, 20]$  is usually sufficient to achieve high performance for most real graphs. Nevertheless, longer random walks may be necessary in e.g., graphs with small  $\mu'$ , especially when the number of labeled nodes is scarce. To deal with such challenges, the ensuing modification of AdaDIF that scales linearly with K is nicely motivated.

# C. Dictionary of diffusions

The present section deals with a modified version of AdaDIF, where the number of parameters (dimension of  $\theta$ ) is restricted to D < K, meaning the "degrees of freedom" of each class-specific distribution are fewer than the number of landing probabilities. Specifically, consider (cf. (4))

$$\mathbf{f}_c(\boldsymbol{\theta}) = \sum_{k=1}^K a_k(\boldsymbol{\theta}) \mathbf{p}_c^{(k)} = \mathbf{P}_c^{(K)} \mathbf{a}(\boldsymbol{\theta})$$

where  $a_k(\boldsymbol{\theta}) := \sum_{d=1}^D \theta_d C_{kd}$ , and  $\mathbf{C} := [\mathbf{c}_1 \cdots \mathbf{c}_D] \in \mathbb{R}^{K \times D}$  is a *dictionary* of D coefficient vectors, the  $i^{th}$  forming the column  $\mathbf{c}_i \in \mathcal{S}^K$ . Since  $\mathbf{a}(\boldsymbol{\theta}) = \mathbf{C}\boldsymbol{\theta}$ , it follows that

$$\mathbf{f}_c(oldsymbol{ heta}) = \mathbf{P}_c^{(K)} \mathbf{C} oldsymbol{ heta} = \sum_{d=1}^D heta_d \mathbf{f}_c^{(d)}$$

where  $\mathbf{f}_c^{(d)} := \sum_{k=1}^K C_{kd} \mathbf{p}_c^{(k)}$  is the  $d^{th}$  diffusion. To find the optimal  $\boldsymbol{\theta}$ , the optimization problem in (12) is

To find the optimal  $\theta$ , the optimization problem in (12) is solved with

$$\mathbf{b}_{c} = -\frac{2}{|\mathcal{L}|} (\mathbf{F}_{c}^{\Delta})^{\mathsf{T}} \mathbf{D}_{\mathcal{L}}^{\dagger} \mathbf{y}_{\mathcal{L}^{c}}$$
(17)

$$\mathbf{A}_{c} = (\mathbf{F}_{c}^{\Delta})^{\mathsf{T}} \mathbf{D}_{\mathcal{L}}^{\dagger} \mathbf{F}_{c}^{\Delta} + \lambda (\mathbf{F}_{c}^{\Delta})^{\mathsf{T}} \mathbf{D}^{-1} \mathbf{L} \mathbf{D}^{-1} \mathbf{F}_{c}^{\Delta}$$
(18)

where  $\mathbf{F}_c^{\Delta} := [\mathbf{f}_c^{(1)} \cdots \mathbf{f}_c^{(D)}]$  effectively replaces  $\mathbf{P}_c^{(K)}$  as the basis of the space on which each  $\mathbf{f}_c$  is constructed. The description of AdaDIF in *dictionary mode* is given as a special case of Algorithm 1, together with the subroutine in Algorithm 3 for memory-efficient generation of  $\mathbf{F}_c^{\Delta}$ .

The motivation behind this dictionary-based variant of AdaDIF is two-fold. First, it leverages the properties of a judiciously selected basis of known diffusions, e.g. by constructing  $\mathbf{C} = [\boldsymbol{\theta}_{\mathrm{PPR}}, \boldsymbol{\theta}_{\mathrm{HK}}, \ldots]$ . In that sense, our approach is related to multi-kernel methods, e.g. [1], although significantly more scalable than the latter. Second, the complexity of AdaDIF in dictionary mode is  $\mathcal{O}(|\mathcal{E}|(K+D))$ , where D can be arbitrarily smaller than K, leading to complexity that is *linear* with respect to both K and  $|\mathcal{E}|$ .

# **Algorithm 1** Adaptive Diffusions

**Input:** Adjacency matrix: W, Labeled nodes:  $\{y_i\}_{i\in\mathcal{L}}$ parameters: Regularization parameter:  $\lambda$ , # of landing probabilities: K, Dictionary mode  $\in \{\text{True}, \text{False}\}$ , Unconstrained  $\in \{\text{True}, \text{False}\}$ **Output:** Predictions:  $\{\hat{y}_i\}_{i\in\mathcal{U}}$ Extract  $\mathcal{Y} = \{$  Set of unique labels in:  $\{y_i\}_{i \in \mathcal{L}}\}$ for  $c \in \mathcal{Y}$  do  $\mathcal{L}_c = \{i \in \mathcal{L} : y_i = c\}$ if Dictionary mode then  $\mathbf{F}_c^{\Delta} = \text{Dictionary} (\mathbf{W}, \mathcal{L}_c, K, \mathbf{C})$ Obtain  $\mathbf{b}_c$  and  $\mathbf{A}_c$  as in (17) and (18) else  $\{\mathbf{P}_c^{(K)}, \tilde{\mathbf{P}}_c^{(K)}\} = \text{LANDPROB}(\mathbf{W}, \mathcal{L}_c, K)$ Obtain  $\mathbf{b}_c$  and  $\mathbf{A}_c$  as in (13) and (15)  $\mathbf{F}_c = \mathbf{P}_c^{(K)}$ end if if Unconstrained then Obtain  $\hat{\boldsymbol{\theta}}_c$  as in (19) else Obtain  $\hat{\boldsymbol{\theta}}_c$  by solving (12)  $\mathbf{f}_c(\hat{\boldsymbol{\theta}}_c) = \mathbf{F}_c\hat{\boldsymbol{\theta}}_c$ Obtain  $\hat{y}_i = \arg\max_{c \in \mathcal{Y}} \left[ \mathbf{f}_c(\hat{\boldsymbol{\theta}}_c) \right]$ ,  $\forall i \in \mathcal{U}$ 

## Algorithm 2 LANDPROB

```
Input: \mathbf{W}, \mathcal{L}_c, K Output: \mathbf{P}_c^{(K)}, \ \tilde{\mathbf{P}}_c^{(K)}
\mathbf{H} = \mathbf{W}\mathbf{D}^{-1}; \ \mathbf{p}_c^{(0)} = \mathbf{v}_c
for k = 1: K+1 do
\mathbf{p}_c^{(k)} = \mathbf{H}\mathbf{p}_c^{(k-1)}
\tilde{\mathbf{p}}_c^{(k)} = \mathbf{p}_c^{(k-1)} - \mathbf{p}_c^{(k)}
end for
```

### Algorithm 3 DICTIONARY

```
Input: W, \mathcal{L}_c, K, C Output: \mathbf{F}_c^{\Delta}

\mathbf{H} = \mathbf{W}\mathbf{D}^{-1}; \mathbf{p}_c^{(0)} = \mathbf{v}_c; \{\mathbf{f}_c^{(d)}\}_{d=1}^D = \mathbf{0}

for k = 1: K do
\mathbf{p}_c^{(k)} = \mathbf{H}\mathbf{p}_c^{(k-1)}

for d = 1: D do
\mathbf{f}_c^{(d)} = \mathbf{f}_c^{(d)} + C_{kd}\mathbf{p}_c^{(k)}
end for
```

# D. Unconstrained diffusions

Thus far, the diffusion coefficients  $\theta$  have been constrained on the K-dimensional probability simplex  $\mathcal{S}^K$ , resulting in sparse solutions  $\hat{\theta}_c$ , as well as  $\mathbf{f}_c(\hat{\theta}_c) \in \mathcal{S}^N$ . The latter also allows each  $\mathbf{f}_c(\theta)$  to be interpreted as a pmf over  $\mathcal{V}$ .

Nevertheless, the simplex constraint imposes a limitation to the model: landing probabilities may only have *non-negative* contribution on the resulting class distribution. Upon relaxing this non-negativity constraint, (12) affords a closed-form solution as

$$\hat{\boldsymbol{\theta}}_c = \mathbf{A}_c^{-1} (\mathbf{b}_c - \lambda^* \mathbf{1}), \qquad \lambda^* = \frac{\mathbf{1}^\mathsf{T} \mathbf{A}_c^{-1} \mathbf{b}_c - 1}{\mathbf{b}^\mathsf{T} \mathbf{A}_c^{-1} \mathbf{b}_c}. \quad (19)$$

Retaining the hyperplane constraint  $\mathbf{1}^T \boldsymbol{\theta} = 1$  prevents the trivial solution  $\boldsymbol{\theta} = \mathbf{0}$ , and forces at least one entry of  $\boldsymbol{\theta}$  to be positive.

### IV. RELATION TO PRIOR WORKS

Following the seminal contribution in [9] that introduced PageRank as a network centrality measure, there has been a vast body of works studying its theoretical properties, computational aspects, as well as applications beyond Web ranking [23], [14]. Most formal approaches to generalize PageRank focus either on the teleportation component (see e.g. [28], [29] as well as [8] for an application to semisupervised classification), or, on the so-termed damping mechanism [12], [4]. Perhaps the most general setting can be found in [4], where a family of functional rankings was introduced by the choice of a parametric damping function that assigns weights to successive steps of a walk initialized according to the teleportation distribution. The per class distributions produced by AdaDIF are in fact members of this family of functional rankings. However, instead of choosing a fixed damping function as in the aforementioned approaches, AdaDIF learns a class-specific and graph-aware damping mechanism. In this sense, AdaDIF undertakes statistical learning in the space of functional rankings, tailored to the underlying semi-supervised classification task. AdaDIF also shares links with SSL methods based on graph signal processing proposed in [33]. Similar to our approach, these graph filter based techniques are parametrized via assigning different weights to a number of consecutive powers of a matrix related to the structure of the graph. Our contribution however, introduces different loss and regularization functions for adapting the diffusions. It also uses the simplex constraint which improves the numerical stability of the involved computations; reduces the search-space of the model (which is beneficial under data scarcity); and makes the model amenable to a rigorous analysis that relates the dimensionality of the feature space to basic graph properties.

# V. EXPERIMENTAL EVALUATION

Our experiments compare the classification accuracy of the novel AdaDIF approach with state-of-the-art alternatives. For the comparisons, we use 6 benchmark labeled graphs whose dimensions and basic attributes are summarized in Table III. All 6 graphs have nodes that belong to multiple classes, while the last 3 are *multilabeled* (each node has *one or more* labels). We evaluate performance of AdaDIF and

the following: i) PPR and HK, which are special cases of AdaDIF as discussed in Section II; ii) Label propagation (LP) [38] iii) Node2vec [16]; iv) Deepwalk [30]; v) Planetoid-G [37]; and, vi) graph convolutional networks (GCNs) [18].

We performed 10-fold cross-validation to select parameters needed by i) - v). For HK, we performed grid search over  $t \in [1, 5, 10, 15]$ . For PPR, we fixed  $\alpha = 0.98$  since it is well documented that  $\alpha$  close to 1 yields reliable performance; see e.g., [25]. HK, PPR and LP were run for 50 steps for convergence to be in effect. For Node2vec, we fixed most parameters to the values suggested in [16], and performed grid search for  $p, q \in [0.25, 1.0, 2.0, 4.0]$ . Since Deepwalk can be seen as Node2vec with p = q = 1.0, we used the Node2vec Python implementation for both. As in [16], [30], we used the embedded node-features to train a supervised logistic regression classifier with  $\ell_2$  regularization. For AdaDIF, we fixed  $\lambda = 15.0$ , while K = 20 was sufficient to attain desirable accuracy; only the values of Boolean variables Unconstrained and Dictionary Mode (see Algorithm 1) were tuned by validation. For the multilabel graphs, we found  $\lambda = 5.0$  and even shorter walks of K = 10to perform well. For the dictionary mode of AdaDIF, we preselected D=10, with the first five columns of C being HK coefficients with parameters  $t \in [5.0, 20.0]$ , and the other five polynomial coefficients  $c_i = k^{\beta}$  with  $\beta \in [2.0, 10]$ .

For multiclass experiments, we evaluated the performance of all algorithms on the three benchmark citation networks, namely Cora, Citeseer, and PubMed. We obtained the labels of an increasing number of nodes by uniformly sampling  $|\mathcal{L}_c|$  nodes from each class, and predicted the labels of the remaining nodes. For each experiment, classification accuracy was measured on the unlabeled nodes in terms of Micro-F1 and Macro-F1 scores; see e.g., [26]. The results were averaged over 20 experiments, with mean and standard deviation reported in Table I. Evidently, AdaDIF achieves state of the art performance for all graphs. For Cora and PubMed, AdaDIF was switched to dictionary mode, while for Citeseer, where the gain in accuracy is more significant, unconstrained diffusions were employed. In the multiclass setting, diffusion-based classifiers (AdaDIF, PPR, and HK) outperformed the embedding-based methods by a small margin, and GCNs by a larger margin. It should be noted however that GCNs were mainly designed to combine the graph with node features. In our "featureless" setting, we used the identity matrix columns as input features, as suggested in [18, Appendix]. The scalabilty of AdaDIF is reflected on the runtime comparisons listed in Fig. 3. All experiments were run on a machine with i5 @3.50 Mhz CPU, and 16GB of RAM. For the compared algorithms we used the implementations provided by the authors.

Finally, Table II presents the results on multilabel graphs, where we compare with Deepwalk and Node2vec, since the rest of the methods are designed for multiclass problems. Since nodes here may have multiple labels, it is challenging

Table I

MICRO F1 AND MACRO F1 SCORES ON MULTICLASS NETWORKS (CLASS-BALANCED SAMPLING)

	Network	Cora			Citeseer			PubMed		
	$ \mathcal{L}_c $	5	10	20	5	10	20	5	10	20
Micro-F1	AdaDIF PPR HK LP Node2vec Deepwalk Planetoid-G GCN	$67.5 \pm 2.2$ $67.1 \pm 2.3$ $67.0 \pm 2.5$ $61.8 \pm 3.5$ $68.9 \pm 1.9$ $68.4 \pm 2.0$ $63.5 \pm 4.7$ $60.1 \pm 3.7$	$71.0 \pm 2.0$ $70.2 \pm 2.1$ $70.5 \pm 2.5$ $66.3 \pm 4.2$ $70.2 \pm 1.6$ $70.0 \pm 1.6$ $65.6 \pm 2.7$ $65.5 \pm 2.5$	$73.2 \pm 1.2 \\ 72.8 \pm 1.5 \\ 72.9 \pm 1.2 \\ 71.0 \pm 2.7 \\ 72.4 \pm 1.2 \\ 72.0 \pm 1.4 \\ 69.0 \pm 1.5 \\ 68.6 \pm 1.9$	$42.3 \pm 4.4 \\ 41.1 \pm 5.2 \\ 40.0 \pm 5.6 \\ 40.7 \pm 2.5 \\ 39.2 \pm 3.7 \\ 38.4 \pm 3.9 \\ 37.8 \pm 4.0 \\ 38.3 \pm 3.2$	$49.5 \pm 3.0 \\ 48.7 \pm 2.5 \\ 48.0 \pm 2.4 \\ 48.0 \pm 3.7 \\ 46.5 \pm 2.4 \\ 45.5 \pm 2.0 \\ 44.9 \pm 3.3 \\ 44.2 \pm 2.2$		$62.0 \pm 6.0$ $63.1 \pm 1.1$ $62.0 \pm 0.6$ $56.2 \pm 11.0$ $61.7 \pm 13.0$ $61.5 \pm 1.3$ $60.7 \pm 2.0$ $60.0 \pm 1.9$	$68.5 \pm 4.5$ $69.5 \pm 3.8$ $68.3 \pm 4.7$ $68.0 \pm 6.1$ $66.4 \pm 4.6$ $65.8 \pm 5.0$ $63.4 \pm 2.3$ $63.6 \pm 2.5$	$74.1 \pm 1.7$ $74.1 \pm 1.8$ $74.0 \pm 1.8$ $69.3 \pm 2.4$ $71.1 \pm 2.4$ $70.5 \pm 2.2$ $68.0 \pm 1.5$ $70.5 \pm 1.5$
Macro-F1	AdaDIF PPR HK LP Node2vec Deepwalk Planetoid-G GCN	$\begin{array}{c} 65.5 \pm 2.5 \\ 65.0 \pm 2.3 \\ 65.0 \pm 2.5 \\ 60.1 \pm 3.2 \\ 62.4 \pm 2.0 \\ 61.8 \pm 2.2 \\ 59.9 \pm 4.5 \\ 53.8 \pm 6.6 \end{array}$	$70.6 \pm 2.2$ $70.0 \pm 2.3$ $70.0 \pm 2.6$ $66.5 \pm 4.1$ $64.7 \pm 1.7$ $64.5 \pm 2.0$ $63.0 \pm 3.0$ $61.9 \pm 2.6$		$36.1 \pm 3.9 \\ 34.7 \pm 5.0 \\ 33.9 \pm 5.4 \\ 34.8 \pm 4.6 \\ 34.6 \pm 2.7 \\ 34.0 \pm 2.5 \\ 33.3 \pm 2.5 \\ 32.8 \pm 2.0$	$44.0 \pm 2.8 \\ 43.5 \pm 2.3 \\ 42.8 \pm 2.2 \\ 41.8 \pm 3.9 \\ 41.6 \pm 1.9 \\ 41.0 \pm 2.0 \\ 40.2 \pm 2.2 \\ 39.1 \pm 1.8$	$48.1 \pm 1.2 \\ 47.6 \pm 0.6 \\ 47.0 \pm 0.6 \\ 51.5 \pm 1.2 \\ 45.3 \pm 1.5 \\ 44.7 \pm 1.8 \\ 43.6 \pm 2.0 \\ 43.0 \pm 1.7$	$\begin{array}{c} 60.4\pm0.6\\ \textbf{61.7}\pm\textbf{0.6}\\ \textbf{60.5}\pm0.6\\ 51.5\pm12.3\\ 59.5\pm1.2\\ 59.3\pm1.2\\ 57.7\pm1.5\\ 54.4\pm4.1 \end{array}$	$67.0 \pm 4.4$ $68.1 \pm 3.6$ $66.8 \pm 4.7$ $66.2 \pm 6.6$ $64.0 \pm 3.8$ $63.8 \pm 4.0$ $61.9 \pm 3.5$ $57.2 \pm 5.2$	$ \begin{array}{c} \textbf{72.6} \pm 1.8 \\ \textbf{72.6} \pm 1.8 \\ \textbf{72.7} \pm 1.8 \\ \textbf{67.8} \pm 2.0 \\ \textbf{72.3} \pm 1.4 \\ \textbf{72.1} \pm 1.3 \\ \textbf{66.1} \pm 1.8 \\ \textbf{60.5} \pm 2.4 \\ \end{array} $

 ${\bf Table~II} \\ {\bf Micro~F1~and~Macro~F1~Scores~of~Various~Algorithms~on~Multilabel~Networks}$ 

	Network	work PPI			BlogCatalog			Wikipedia		
	$ \mathcal{L} / \mathcal{V} $	10%	20%	30%	10%	20%	30%	10%	20%	30%
Micro-F1	AdaDIF PPR HK Node2vec Deepwalk	$15.4 \pm 0.5$ $13.8 \pm 0.5$ $14.5 \pm 0.5$ $16.5 \pm 0.6$ $16.0 \pm 0.6$	$17.9 \pm 0.7$ $15.8 \pm 0.6$ $16.7 \pm 0.6$ $18.2 \pm 0.3$ $17.9 \pm 0.5$	$19.2 \pm 0.6$ $17.0 \pm 0.4$ $18.1 \pm 0.5$ $19.1 \pm 0.3$ $18.8 \pm 0.4$	$31.5 \pm 0.6$ $21.1 \pm 0.8$ $22.2 \pm 1.0$ $35.0 \pm 0.3$ $34.2 \pm 0.4$	$34.4 \pm 0.5$ $23.6 \pm 0.6$ $24.7 \pm 0.7$ $36.3 \pm 0.3$ $35.7 \pm 0.3$	$36.3 \pm 0.4$ $25.2 \pm 0.6$ $26.6 \pm 0.7$ $37.2 \pm 0.2$ $36.4 \pm 0.4$	$28.2 \pm 0.9$ $10.5 \pm 1.5$ $9.3 \pm 1.4$ $42.3 \pm 0.9$ $41.0 \pm 0.8$	$30.0 \pm 0.5$ $8.1 \pm 0.7$ $7.3 \pm 0.7$ $44.0 \pm 0.6$ $43.5 \pm 0.5$	$31.2 \pm 0.7$ $7.2 \pm 0.5$ $6.0 \pm 0.7$ $45.1 \pm 0.4$ $44.1 \pm 0.5$
Macro-F1	AdaDIF PPR HK Node2vec Deepwalk	$13.4 \pm 0.6 \\ 12.9 \pm 0.4 \\ 13.4 \pm 0.6 \\ 13.1 \pm 0.6 \\ 12.7 \pm 0.7$	$15.4 \pm 0.7 \\ 14.7 \pm 0.5 \\ 15.4 \pm 0.5 \\ 15.2 \pm 0.5 \\ 15.1 \pm 0.6$	$16.5 \pm 0.7 \\ 15.8 \pm 0.4 \\ 16.5 \pm 0.4 \\ 16.0 \pm 0.5 \\ 16.0 \pm 0.5$	$\begin{array}{c} \textbf{23.0} \pm \textbf{0.6} \\ 17.3 \pm 0.5 \\ 18.4 \pm 0.6 \\ 16.8 \pm 0.5 \\ 16.6 \pm 0.5 \end{array}$	$\begin{array}{c} \textbf{25.3} \pm \textbf{0.4} \\ 19.5 \pm 0.4 \\ 20.7 \pm 0.4 \\ 19.0 \pm 0.3 \\ 18.7 \pm 0.5 \end{array}$	$\begin{array}{c} \textbf{27.0} \pm \textbf{0.4} \\ 20.8 \pm 0.3 \\ 22.3 \pm 0.4 \\ 20.1 \pm 0.4 \\ 19.6 \pm 0.4 \end{array}$	$7.7 \pm 0.3$ $4.4 \pm 0.3$ $4.2 \pm 0.4$ $7.6 \pm 0.3$ $7.3 \pm 0.3$	$8.3 \pm 0.3$ $3.8 \pm 0.6$ $3.7 \pm 0.5$ $8.2 \pm 0.3$ $8.1 \pm 0.2$	$9.0 \pm 0.2 \\ 3.6 \pm 0.2 \\ 3.5 \pm 0.2 \\ 8.5 \pm 0.3 \\ 8.2 \pm 0.2$

Table III
NETWORK CHARACTERISTICS

Network	$ \mathcal{V} $	$ \mathcal{E} $	$ \mathcal{Y} $	Multilabel	
Citeseer	3,233	9,464	6	No	
Cora	2,708	10,858	7	No	
PubMed	19,717	88,676	3	No	
PPI (H. Sapiens)	3,890	76,584	50	Yes	
Wikipedia	4,733	184,182	40	Yes	
BlogCatalog	10,312	333,983	39	Yes	

to find a subset of nodes that has equal number of labels for each class. In fact, such a set may not even exist for a given size. Therefore, for multilabel graphs we simply draw nodes uniformly at random and observe their labels. Also, since these graphs have a large number of classes, we increased the number of training samples. Similar to [16] and [30], during evaluation the number of labels per sampled node is known, and check how many of them are in the top predictions. First, we observe that AdaDIF markedly outperforms PPR and HK across graphs and metrics. Furthermore, for the PPI and BlogCatalog graphs the Micro-F1 score of AdaDIF comes close to that of the much heavier state-of-the-art Node2vec. Finally, AdaDIF outperforms the competing alternatives in terms of Macro-F1 score.

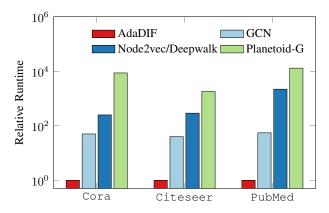


Figure 3. Relative runtime for multiclass networks.

## VI. CONCLUSIONS

We introduced a principled, data-efficient approach to learning class-specific diffusion functions tailored to the underlying network topology. Experiments on real networks confirm that adapting the diffusion function to the given graph and observed labels, significantly improves the performance over fixed diffusions; reaching—and many times surpassing—the classification accuracy of state-of-the-art competing methods while being orders of magnitude faster.

#### REFERENCES

- [1] A. Argyriou, M. Herbster, and M. Pontil, "Combining graph laplacians for semi–supervised learning," in *Proc. Advances in Neural Information Processing Systems*, Vancouver, Can., 2006, pp. 67–74.
- [2] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," in *Proc. Advances in Neural Information Process*ing Systems, Barcelona, Spain, 2016, pp. 1993–2001.
- [3] K. Avrachenkov, A. Mishenin, P. Gonçalves, and M. Sokol, "Generalized optimization framework for graph-based semisupervised learning," *Proc. SIAM Intl. Conf. on Data Mining*, Anaheim, CA, 2012, pp. 966–974.
- [4] R. Baeza-Yates, P. Boldi, and C. Castillo, "Generic damping functions for propagating importance in link-based ranking," *Internet Math.*, vol. 3, no. 4, pp. 445–478, 2006.
- [5] M. Belkin, P. Niyogi, and V. Sindhwani, "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *J. Mach. Learn. Res.*, no. 7, Nov, 2006, pp. 2399–2434.
- [6] Y. Bengio, O. Delalleau, and N. Le Roux, "Label propagation and quadratic criterion," in *Semi-Supervised Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [7] D. Berberidis, A. N. Nikolakopoulos, and G. B. Giannakis, "Adaptive Diffusions for Scalable Learning over Graphs," arXiv preprint arXiv:1804.02081.
- [8] D. Berberidis, A. N. Nikolakopoulos, and G. B. Giannakis, "Random walks with restarts for graph-based classification: Teleportation tuning and sampling design," in *Proc. of IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, Calgary, Can., April 2018.
- [9] S. Brin and L. Page, "Reprint of: The anatomy of a large-scale hypertextual web search engine," *Comput. Netw.*, vol. 56, no. 18, pp. 3825–3833, 2012.
- [10] E. Buchnik and E. Cohen, "Bootstrapped graph diffusions: Exposing the power of nonlinearity," arXiv preprint arXiv:1703.02618, 2017.
- [11] O. Chapelle, B. Schölkopf, and A. Zien, Semi-Supervised Learning. Cambridge, MA, USA: MIT Press, 2006.
- [12] P. G. Constantine and D. F. Gleich, "Random alpha pagerank," Internet Math., vol. 6, no. 2, pp. 189–236, 2009.
- [13] F. Chung, "The heat kernel as the pagerank of a graph," Proc. Natl. Acad. Sci., vol. 104, no. 50, pp. 19735–19740, 2007.
- [14] D. F. Gleich, "Pagerank beyond the web," SIAM Rev., vol. 57, no. 3, pp. 321–363, 2015.
- [15] J. Gorski, F. Pfeuffer, and K. Klamroth, "Biconvex sets and optimization with biconvex functions: a survey and extensions," *Math. Methods of Oper. Res.*, vol. 66, no. 3, pp. 373–407, Dec. 2007.
- [16] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *Proc. of ACM SIGKDD Int. Conf.* on Knowledge Discovery and Data Mining, San Francisco, CA, 2016, pp. 855–864.
- [17] T. Joachims, "Transductive learning via spectral graph partitioning," *Proc. of Intl. Conf. on Machine Learn.*, Washington DC, 2003, pp. 290–297.
- [18] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," arXiv preprint arXiv:1609.02907, 2016.
- [19] K. Kloster and D. F. Gleich, "Heat kernel based community detection," in *Proc. of ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, New York, NY, 2014, pp. 1386– 1395.

- [20] I. M. Kloumann, J. Ugander, and J. Kleinberg, "Block models and personalized pagerank," *Proc. Natl. Acad. Sci.*, vol. 114, no. 1, pp. 33–38, 2017.
- [21] R. I. Kondor and J. Lafferty, "Diffusion kernels on graphs and other discrete input spaces," in *Proc. of Int. Conf. on Machine Learning*, Syndey, Australia, 2002, pp. 315–322.
- [22] B. Kveton, M. Valko, A. Rahimi, and L. Huang, "Semi-supervised learning with max-margin graph cuts," in *Proc. of. Int. Conf. on Artificial Intelligence and Statistics*, Sardinia, Italy, 2010, pp. 421–428.
- [23] A. N. Langville and C. D. Meyer, "Deeper inside pagerank," Internet Math., vol. 1, no. 3, pp. 335–380, 2004.
- [24] D. A. Levin and Y. Peres, Markov Chains and Mixing Times. New York, NY, USA: Amer. Math. Soc., 2017.
- [25] F. Lin and W. W. Cohen, "Semi-supervised classification of network data using very few labels," in *Proc. of Int. Conf. on Advances in Social Network Analysis and Mining*, Odense, Denmark, 2010, pp. 192–199.
- [26] C. D. Manning, P. Raghavan, and H. Schutze, *Introduction to Information Retrieval*. Cambridge, MA: Cambridge University Press, 2008.
- [27] E. Merkurjev, A. L. Bertozzi, and F. Chung, "A semisupervised heat kernel pagerank MBO algorithm for data classification," Univ. of California Los Angeles, Los Angeles, US, Tech. Rep., 2016.
- [28] A. N. Nikolakopoulos and J. D. Garofalakis, "Ncdawarerank: A novel ranking method that exploits the decomposable structure of the web," *Proc. ACM Intl. Conf. on Web Search and Data Mining*, Rome, Italy, 2013, pp. 143–152.
- [29] A. N. Nikolakopoulos, A. Korba, and J. D. Garofalakis, "Random surfing on multipartite graphs," in *Proc. of IEEE Int. Conf. on Big Data*, Washington DC, Dec. 2016, pp. 736–745.
- [30] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," *Proc. ACM SIGKDD Intl. Conf. on Knowl. Disc. and Data Mining*, New York, NY, 2014, pp. 701–710.
- [31] A. T. Puig, A. Wiesel, G. Fleury, and A. O. Hero, "Multidimensional shrinkage-thresholding operator and group lasso penalties," *IEEE Signal Process. Lett.*, vol. 18, no. 6, pp. 363–366, 2011.
- [32] N. Rosenfeld and A. Globerson, "Semi-supervised learning with competitive infection models," arXiv preprint arXiv:1703.06426, 2017.
- [33] A. Sandryhaila and J. M. F. Moura, "Discrete signal processing on graphs," *IEEE Trans. Signal Process.*, vol. 61, no. 7, pp. 1644–1656, April 2013.
- [34] P. P. Talukdar and K. Crammer, "New regularized algorithms for transductive learning," in *Proc. of Joint Eur. Conf. on Machine Learning and Knowledge Discovery in Databases*, 2009, pp. 442–457.
- [35] J. Ugander and L. Backstrom, "Balanced label propagation for partitioning massive graphs," in *Proc. of ACM Int. Conf. on Web Search and Data Mining*, Rome, Italy, 2013, pp. 507–516.
- [36] X.-M. Wu, Z. Li, A. M. So, J. Wright, and S.-F. Chang, "Learning with partially absorbing random walks," *Proc. Adv. in Neural Inform. Proc. Systems*, Lake Tahoe, CA, Dec. 2012, pp. 3077–3085.
- [37] Z. Yang, W. W. Cohen, and R. Salakhutdinov, "Revisiting semisupervised learning with graph embeddings," arXiv preprint arXiv:1603.08861, 2016.
- [38] X. Zhu, Z. Ghahramani, and J. Lafferty, "Semi-supervised learning using Gaussian fields and harmonic functions," in *Proc. of Int. Conf. on Machine Learning*, Washington DC, Aug. 2003.