

---

# Generalized Matrix Means for Semi-Supervised Learning with Multilayer Graphs

---

Pedro Mercado<sup>1</sup>, Francesco Tudisco<sup>2</sup> and Matthias Hein<sup>1</sup>

<sup>1</sup>University of Tübingen, Germany

<sup>2</sup>Gran Sasso Science Institute, Italy

## Abstract

We study the task of semi-supervised learning on multilayer graphs by taking into account both labeled and unlabeled observations together with the information encoded by each individual graph layer. We propose a regularizer based on the generalized matrix mean, which is a one-parameter family of matrix means that includes the arithmetic, geometric and harmonic means as particular cases. We analyze it in expectation under a Multilayer Stochastic Block Model and verify numerically that it outperforms state of the art methods. Moreover, we introduce a matrix-free numerical scheme based on contour integral quadratures and Krylov subspace solvers that scales to large sparse multilayer graphs.

## 1 Introduction

The task of graph-based Semi-Supervised Learning (SSL) is to build a classifier that takes into account both labeled and unlabeled observations, together with the information encoded by a given graph[4, 27]. A common and successful approach is to take a suitable loss function on the labeled nodes and a regularizer which provides information encoded by the graph [2, 15, 30, 32, 35]. Whereas this task is well studied, traditionally these methods assume that the graph is composed by interactions of one single kind, i.e. only one graph is available.

For the case where multiple graphs, or equivalently, multiple layers are available, the challenge is to boost the classification performance by merging the information encoded in each graph. The arguably most popular approach for this task consists of finding some form of convex combination of graph matrices, where more informative graphs receive a larger weight [1, 13, 14, 23, 28, 29, 31, 33].

Note that a convex combination of graph matrices can be seen as a weighted arithmetic mean of graph matrices. In the context of multilayer graph clustering, previous studies [19–21] have shown that weighted arithmetic means are suboptimal under certain benchmark generative graph models, whereas other matrix means, such as the geometric [20] and harmonic means [19], are able to discover clustering structures that the arithmetic means overlook.

In this paper we study the task of semi-supervised learning with multilayer graphs with a novel regularizer based on the power mean Laplacian. The power mean Laplacian is a one-parameter family of Laplacian matrix means that includes as special cases the arithmetic, geometric and harmonic mean of Laplacian matrices. We show that in expectation under a Multilayer Stochastic Block Model, our approach provably correctly classifies unlabeled nodes in settings where state of the art approaches fail. In particular, a limit case of our method is provably robust against noise, yielding good classification performance as long as one layer is informative and remaining layers are potentially just noise. We verify the analysis in expectation with extensive experiments with random graphs, showing that our approach compares favorably with state of the art methods, yielding a good classification performance on several relevant settings where state of the art approaches fail.

name	minimum	harmonic mean	geometric mean	arithmetic mean	maximum
$p$	$p \rightarrow -\infty$	$p = -1$	$p \rightarrow 0$	$p = 1$	$p \rightarrow \infty$
$m_p(a, b)$	$\min\{a, b\}$	$2\left(\frac{1}{a} + \frac{1}{b}\right)^{-1}$	$\sqrt{ab}$	$(a + b)/2$	$\max\{a, b\}$

Table 1: Particular cases of scalar power means

Moreover, our approach scales to large datasets: even though the computation of the power mean Laplacian is in general prohibitive for large graphs, we present a matrix-free numerical scheme based on integral quadratures methods and Krylov subspace solvers which allows us to apply the power mean Laplacian regularizer to large sparse graphs. Finally, we perform numerical experiments on real world datasets and verify that our approach is competitive to state of the art approaches.

## 2 The Power Mean Laplacian

In this section we introduce our multilayer graph regularizer based on the power mean Laplacian. We define a multilayer graph  $\mathbb{G}$  with  $T$  layers as the set  $\mathbb{G} = \{G^{(1)}, \dots, G^{(T)}\}$ , with each graph layer defined as  $G^{(t)} = (V, W^{(t)})$ , where  $V = \{v_1, \dots, v_n\}$  is the node set and  $W^{(t)} \in \mathbb{R}_+^{n \times n}$  is the corresponding adjacency matrix, which we assume symmetric and nonnegative. We further denote the layers' normalized Laplacians as  $L_{\text{sym}}^{(t)} = I - (D^{(t)})^{-1/2}W^{(t)}(D^{(t)})^{-1/2}$ , where  $D^{(t)}$  is the degree diagonal matrix with  $(D^{(t)})_{ii} = \sum_{j=1}^n W_{ij}^{(t)}$ .

The scalar power mean is a one-parameter family of scalar means defined as

$$m_p(x_1, \dots, x_T) = \left(\frac{1}{T} \sum_{i=1}^T x_i^p\right)^{1/p}$$

where  $x_1, \dots, x_T$  are nonnegative scalars and  $p$  is a real parameter. Particular choices of  $p$  yield specific means such as the arithmetic, geometric and harmonic means, as illustrated in Table 1.

The **Power Mean Laplacian**, introduced in [19], is a matrix extension of the scalar power mean applied to the Laplacians of a multilayer graph and proposed as a more robust way to blend the information encoded across the layers. It is defined as

$$L_p = \left(\frac{1}{T} \sum_{i=1}^T (L_{\text{sym}}^{(i)})^p\right)^{1/p}$$

where  $A^{1/p}$  is the unique positive definite solution of the matrix equation  $X^p = A$ . For the case  $p \leq 0$  a small diagonal shift  $\varepsilon > 0$  is added to each Laplacian, i.e. we replace  $L_{\text{sym}}^{(i)}$  with  $L_{\text{sym}}^{(i)} + \varepsilon$ , to ensure that  $L_p$  is well defined as suggested in [3]. In what follows all the proofs hold for an arbitrary shift. Following [19], we set  $\varepsilon = \log_{10}(1 + |p|) + 10^{-6}$  for  $p \leq 0$  in the numerical experiments.

## 3 Multilayer Semi-Supervised Learning with the Power Mean Laplacian

In this paper we consider the following optimization problem for the task of semi-supervised learning in multilayer graphs: Given  $k$  classes  $r = 1, \dots, k$  and membership vectors  $Y^{(r)} \in \mathbb{R}^n$  defined by  $Y_i^{(r)} = 1$  if node  $v_i$  belongs to class  $r$  and  $Y_i^{(r)} = 0$  otherwise, we let

$$f^{(r)} = \arg \min_{f \in \mathbb{R}^n} \|f - Y^{(r)}\|^2 + \lambda f^T L_p f. \quad (1)$$

The final class assignment for an unlabeled node  $v_i$  is  $y_i = \arg \max\{f_i^{(1)}, \dots, f_i^{(k)}\}$ . Note that the solution  $f$  of (1), for a particular class  $r$ , is such that  $(I + \lambda L_p)f = Y^{(r)}$ . Equation (1) has two terms: the first term is a loss function based on the labeled nodes whereas the second term is a regularization term based on the power mean Laplacian  $L_p$ , which accounts for the multilayer graph structure. It is worth noting that the Local-Global approach of [32] is a particular case of our approach when only one layer ( $T = 1$ ) is considered. Moreover, note that when  $p = 1$  we obtain a regularizer term based on the arithmetic mean of Laplacians  $L_1 = \frac{1}{T} \sum_{i=1}^T L_{\text{sym}}^{(i)}$ . In the following section we analyze our proposed approach (1) under the Multilayer Stochastic Block Model.

## 4 Multilayer Stochastic Block Model

In this section we provide an analysis of semi-supervised learning for multilayer graphs with the power mean Laplacian as a regularizer under the Multilayer Stochastic Block Model (**MSBM**). The MSBM is a generative model for graphs showing certain prescribed clusters/classes structures via a set of membership parameters  $p_{\text{in}}^{(t)}$  and  $p_{\text{out}}^{(t)}$ ,  $t = 1, \dots, T$ . These parameters designate the edge probabilities: given nodes  $v_i$  and  $v_j$  the probability of observing an edge between them on layer  $t$  is  $p_{\text{in}}^{(t)}$  (resp.  $p_{\text{out}}^{(t)}$ ), if  $v_i$  and  $v_j$  belong to the same (resp. different) cluster/class. Note that, unlike the Labeled Stochastic Block Model [11], the MSBM allows multiple edges between the same pairs of nodes across the layers. For SSL with one layer under the SBM we refer the reader to [12, 22, 26].

We present an analysis in expectation. We consider  $k$  clusters/classes  $\mathcal{C}_1, \dots, \mathcal{C}_k$  of equal size  $|\mathcal{C}| = n/k$ . We denote with calligraphic letters the layers of a multilayer graph in expectation  $E(\mathbb{G}) = \{E(G^{(1)}, \dots, E(G^{(T)})\}$ , i.e.  $\mathcal{W}^{(t)}$  is the expected adjacency matrix of the  $t^{\text{th}}$ -layer. We assume that our multilayer graphs are non-weighted, i.e. edges are zero or one, and hence we have  $\mathcal{W}_{ij}^{(t)} = p_{\text{in}}^{(t)}$ , (resp.  $\mathcal{W}_{ij}^{(t)} = p_{\text{out}}^{(t)}$ ) for nodes  $v_i, v_j$  belonging to the same (resp. different) cluster/class.

In order to grasp how different methods classify the nodes in multilayer graphs following the MSBM we analyze two different settings. In the first setting (Section 4.1) all layers have the same class structure and we study the conditions for different regularizers  $L_p$  to correctly predict class labels. We further show that our approach is robust against the presence of noise layers, in the sense that it achieves a small classification error when at least one layer is informative and the remaining layers are potentially just noise. In this setting we distinguish the case where each class has the same amount of initial labels and the case where different classes have different number of labels. In the second setting (Section 4.2) we consider the case where each layer taken alone would lead to a large classification error whereas considering all the layers together can lead to a small classification error.

### 4.1 Complementary Information Layers

A common assumption in multilayer semi-supervised learning is that at least one layer encodes relevant information in the label prediction task. The next theorem discusses the classification error of the expected power mean Laplacian regularizer in this setting.

**Theorem 1.** *Let  $E(\mathbb{G})$  be the expected multilayer graph with  $T$  layers following the multilayer SBM with  $k$  classes  $\mathcal{C}_1, \dots, \mathcal{C}_k$  of equal size and parameters  $(p_{\text{in}}^{(t)}, p_{\text{out}}^{(t)})_{t=1}^T$ . Assume the same number of labeled nodes are available per class. Then, the solution of (1) yields zero test error if and only if*

$$m_p(\boldsymbol{\rho}_\epsilon) < 1 + \epsilon, \quad (2)$$

where  $(\boldsymbol{\rho}_\epsilon)_t = 1 - (p_{\text{in}}^{(t)} - p_{\text{out}}^{(t)}) / (p_{\text{in}}^{(t)} + (k-1)p_{\text{out}}^{(t)}) + \epsilon$ , and  $t = 1, \dots, T$ .

This theorem shows that the power mean Laplacian regularizer allows to correctly classify the nodes if  $p$  is such that condition (2) holds. In order to better understand how this condition changes when  $p$  varies, we analyze in the next corollary the limit cases  $p \rightarrow \pm\infty$ .

**Corollary 1.** *Let  $E(\mathbb{G})$  be an expected multilayer graph as in Theorem 1. Then,*

- For  $p \rightarrow \infty$ , the test error is zero if and only if  $p_{\text{out}}^{(t)} < p_{\text{in}}^{(t)}$  for all  $t = 1, \dots, T$ .
- For  $p \rightarrow -\infty$ , the test error is zero if and only there exists a  $t \in \{1, \dots, T\}$  such that  $p_{\text{out}}^{(t)} < p_{\text{in}}^{(t)}$ .

This corollary implies that the limit case  $p \rightarrow \infty$  requires that *all layers* convey information regarding the clustering/class structure of the multilayer graph, whereas the case  $p \rightarrow -\infty$  requires that *at least one layer* encodes clustering/class information, and hence it is clear that conditions for the limit  $p \rightarrow -\infty$  are less restrictive than the conditions for the limit case  $p \rightarrow \infty$ . The next Corollary shows that the smaller the power parameter  $p$  is, the less restrictive are the conditions to yield a zero test error.

**Corollary 2.** *Let  $E(\mathbb{G})$  be an expected multilayer graph as in Theorem 1. Let  $p \leq q$ . If  $\mathcal{L}_q$  yields zero test error, then  $\mathcal{L}_p$  yields a zero test error.*

The previous results show the effectivity of the power mean Laplacian regularizer in expectation. We now present a numerical evaluation based on Theorem 1 and Corollaries 1 and 2 on random

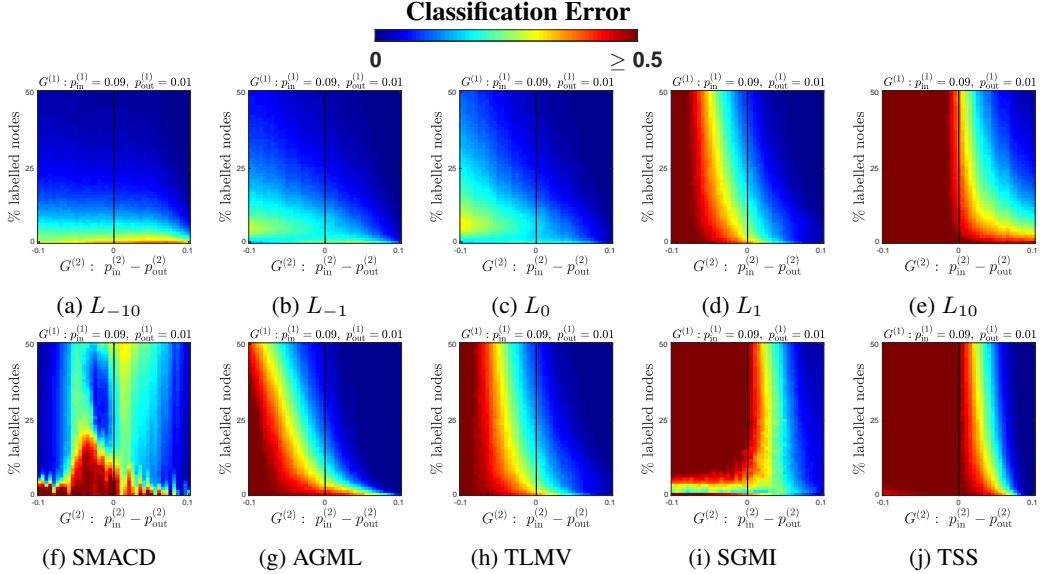
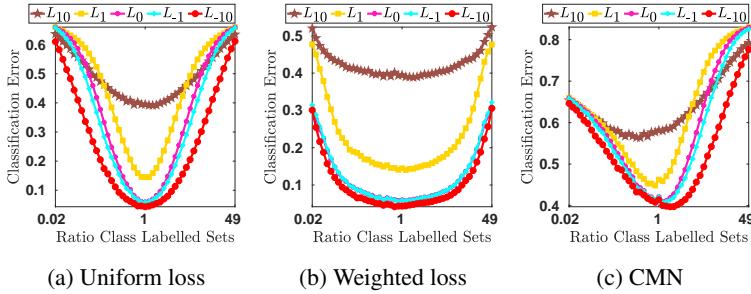


Figure 1: Average classification error under the Stochastic Block Model computed from 100 runs. **Top Row:** Particular cases with the power mean Laplacian. **Bottom Row:** State of the art models.

graphs sampled from the SBM. The corresponding results are presented in Fig. 1 for classification with regularizers  $L_{-10}, L_{-1}, L_0, L_1, L_{10}$  and  $\lambda = 1$ . We first describe the setting we consider: we generate random multilayer graphs with two layers ( $T = 2$ ) and two classes ( $k = 2$ ) each composed by 100 nodes ( $|\mathcal{C}| = 100$ ). For each parameter configuration  $(p_{in}^{(1)}, p_{out}^{(1)}, p_{in}^{(2)}, p_{out}^{(2)})$  we generate 10 random multilayer graphs and 10 random samples of labeled nodes, yielding a total of 100 runs per parameter configuration, and report the average test error. Our goal is to evaluate the classification performance under different SBM parameters and different amounts of labeled nodes. To this end, we fix the first layer  $G^{(1)}$  to be informative of the class structure ( $p_{in}^{(1)} - p_{out}^{(1)} = 0.08$ ), i.e. one can achieve a low classification error by taking this layer alone, provided sufficiently many labeled nodes are given. The second layer will go from non-informative (noisy) configurations ( $p_{in}^{(2)} < p_{out}^{(2)}$ , left half of  $x$ -axis) to informative configurations ( $p_{in}^{(2)} > p_{out}^{(2)}$ , right half of  $x$ -axis), with  $p_{in}^{(t)} + p_{out}^{(t)} = 0.1$  for both layers. Moreover, we consider different amounts of labeled nodes: going from 1% to 50% ( $y$ -axis). The corresponding results are presented in Figs. 1a, 1b, 1c, 1d, and 1e.

In general one can expect a low classification error when both layers  $G^{(1)}$  and  $G^{(2)}$  are informative (right half of  $x$ -axis). We can see that this is the case for all power mean Laplacian regularizers here considered (see top row of Fig. 1). In particular, we can see in Fig. 1e that  $L_{10}$  performs well only when **both** layers are informative and completely fails when the second layer is not informative, regardless of the amount of labeled nodes. On the other side we can see in Fig. 1a that  $L_{-10}$  achieves in general a low classification error, regardless of the configuration of the second layer  $G^{(2)}$ , i.e. when  $G^{(1)}$  or  $G^{(2)}$  are informative. Moreover, we can see that overall the areas with low classification error (dark blue) increase when the parameter  $p$  decreases, verifying the result from Corollary 2. In the bottom row of Fig. 1 we present the performance of state of the art methods. We can observe that most of them present a classification performance that resembles the one of the power mean Laplacian regularizer  $L_1$ . In general their classification performance drops when the level of noise increases, i.e. for non-informative configurations of the second layer  $G^{(2)}$ , and they are outperformed by the power mean Laplacian regularizer for small values of  $p$ .

**Unbalanced Class Proportion on Labeled Datasets.** In the previous analysis we assumed that we had the same amount of labeled nodes per class. We consider now the case where the number of labeled nodes per class is different. This setting was considered in [35], where the goal was to overcome unbalanced class proportions in labeled nodes. To this end, they propose a Class Mass Normalization (CMN) strategy, whose performance was also tested in [34]. In the following result we show that, provided the ground truth classes have the same size, different amounts of labeled nodes per class affect the conditions in expectation for zero classification error of (1). For simplicity, we consider here only the case of two classes.



**Figure 2:** Different class weighted loss strategies. Left to right: uniform loss, weighted loss, and Class Mass Normalization.

**Theorem 2.** Let  $E(\mathbb{G})$  be the expected multilayer graph with  $T$  layers following the multilayer SBM with two classes  $\mathcal{C}_1, \mathcal{C}_2$  of equal size and parameters  $\left(p_{\text{in}}^{(t)}, p_{\text{out}}^{(t)}\right)_{t=1}^T$ . Assume  $n_1, n_2$  nodes from  $\mathcal{C}_1, \mathcal{C}_2$  are labeled, respectively. Let  $\lambda = 1$ . Then (1) yields zero test error if

$$m_p(\rho_\epsilon) < \min \left\{ \frac{n_1}{n_2}, \frac{n_2}{n_1} \right\} \quad (3)$$

where  $(\rho_\epsilon)_t = 1 - (p_{\text{in}}^{(t)} - p_{\text{out}}^{(t)}) / (p_{\text{in}}^{(t)} + (k-1)p_{\text{out}}^{(t)}) + \epsilon$ , and  $t = 1, \dots, T$ .

Observe that Theorem 2 provides only a sufficient condition. A necessary and sufficient condition for zero test error in terms of  $p, n_1$  and  $n_2$  is given in the supplementary material.

A different objective function can be employed for the case of classes with different number of labels per class. Let  $C$  be the diagonal matrix defined by  $C_{ii} = n/n_r$ , if node  $v_i$  has been labeled to belong to class  $\mathcal{C}_r$ . Consider the following modification of (1)

$$\arg \min_{f \in \mathbb{R}^n} \|f - CY\|^2 + \lambda f^T L_p f \quad (4)$$

The next Theorem shows that using (4) in place of (1) allows us to retrieve the same condition of Theorem 1 for zero test error in expectation in the setting where the number of labeled nodes per class are not equal.

**Theorem 3.** Let  $E(\mathbb{G})$  be the expected multilayer graph with  $T$  layers following the multilayer SBM  $k$  classes  $\mathcal{C}_1, \dots, \mathcal{C}_k$  of equal size and parameters  $\left(p_{\text{in}}^{(t)}, p_{\text{out}}^{(t)}\right)_{t=1}^T$ . Let  $n_1, \dots, n_k$  be the number of labeled nodes per class. Let  $C \in \mathbb{R}^{n \times n}$  be a diagonal matrix with  $C_{ii} = n/n_r$  for  $v_i \in \mathcal{C}_r$ . The solution to (4) yields a zero test classification error if and only if

$$m_p(\rho_\epsilon) < 1 + \epsilon, \quad (5)$$

where  $(\rho_\epsilon)_t = 1 - (p_{\text{in}}^{(t)} - p_{\text{out}}^{(t)}) / (p_{\text{in}}^{(t)} + (k-1)p_{\text{out}}^{(t)}) + \epsilon$ , and  $t = 1, \dots, T$ .

In Figs. 2a, 2b, and 2c. we present a numerical experiment with random graphs of our analysis in expectation. We consider the following setting: we generate multilayer graphs with two layers ( $T = 2$ ) and two classes ( $k = 2$ ) each composed by 100 nodes ( $|\mathcal{C}| = 100$ ). We fix  $p_{\text{in}}^{(1)} - p_{\text{out}}^{(1)} = 0.08$  and  $p_{\text{in}}^{(2)} - p_{\text{out}}^{(2)} = 0$ , with  $p_{\text{in}}^{(t)} + p_{\text{out}}^{(t)} = 0.1$  for both layers. We fix the total amount of labeled nodes to be  $n_1 + n_2 = 50$  and let  $n_1, n_2 = 1, \dots, 49$ . For each setting we generate 10 multilayer graphs and 10 sets of labeled nodes, yielding a total of 100 runs per setting, and report the average test classification error. In Fig. 2a we can see the performance of the power mean Laplacian regularizer without modifications. We can observe how different proportions of labeled nodes per class affect the performance. In Fig. 2b, we present the performance of the modified approach (4) and observe that it yields a better performance against different class label proportions. Finally in Fig. 2c we present the performance based on Class Mass Normalization<sup>1</sup>, where we can see that its effect is slightly skewed to one class and its overall performance is larger than the proposed approach.

## 4.2 Information-Independent Layers

In the previous section we considered the case where at least one layer had enough information to correctly estimate node class labels. In this section we now consider the case where single layers

<sup>1</sup>We follow the authors' implementation: [http://pages.cs.wisc.edu/~jerryzhu/pub/harmonic\\_function.m](http://pages.cs.wisc.edu/~jerryzhu/pub/harmonic_function.m)

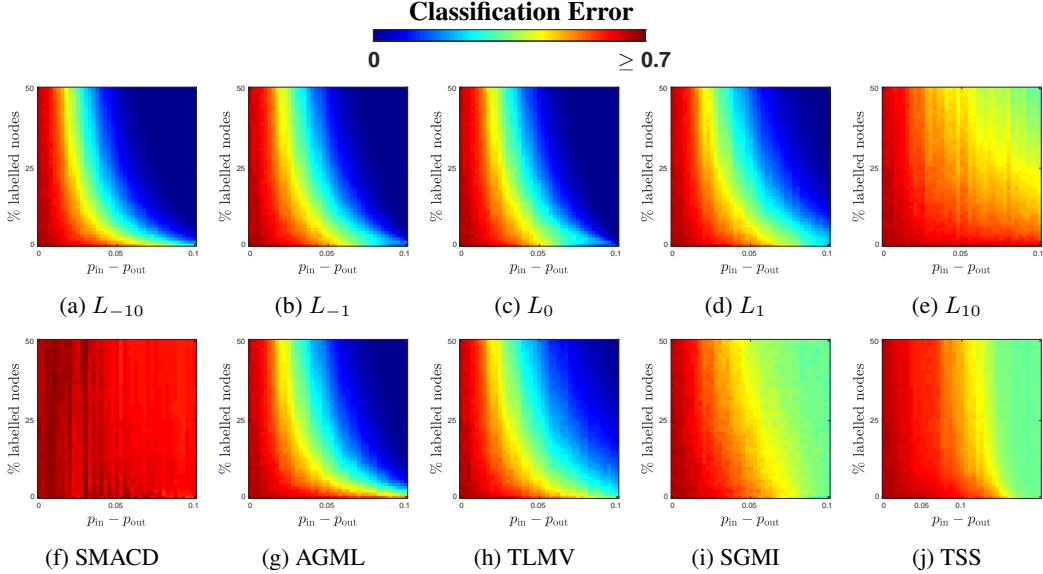


Figure 4: Average test error under the SBM.Multilayer graph with 3 layers and 3 classes. **Top Row:** Particular cases with the power mean Laplacian. **Bottom Row:** State of the art models.

taken alone obtain a large classification error, whereas when all the layers are taken together it is possible to obtain a good classification performance. For this setting we consider multilayer graphs with 3 layers ( $T = 3$ ) and three classes ( $k = 3$ )  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$ , each composed by 100 nodes ( $|\mathcal{C}| = 100$ ) with the following expected adjacency matrix per layer:

$$\mathcal{W}_{i,j}^{(t)} = \begin{cases} p_{in}, & v_i, v_j \in \mathcal{C}_t \text{ or } v_i, v_j \in \overline{\mathcal{C}_t} \\ p_{out}, & \text{else} \end{cases} \quad (6)$$

for  $t = 1, 2, 3$ , i.e. layer  $G^{(t)}$  is informative of class  $\mathcal{C}_t$  but not of the remaining classes, and hence any classification method using one single layer will provide a poor classification performance. In Fig. 4 we present numerical experiments: for each parameter setting  $(p_{in}, p_{out})$  we generate 5 multilayer graphs together with 5 samples of labeled nodes yielding a total of 25 runs per setting, and report the average test classification error. Also in this case we observe that the power mean Laplacian regularizer does identify the global class structure and that it leverages the information provided by labeled nodes, particularly for smaller values of  $p$ . On the other hand, this is not the case for all other state of the art methods. In fact, we can see that SGMI and TSS performs similarly to  $L_{10}$  which has the largest classification error. Moreover, we can see that AGML and TLMV perform similarly to the arithmetic mean of Laplacians  $L_1$ , which in turn is outperformed by the power mean Laplacian regularizer  $L_{-10}$ . Please see the supplementary material for a more detailed comparison.

## 5 A Scalable Matrix-free Numerical Method for the System $(I + \lambda L_p)f = Y$

In this section we introduce a matrix-free method for the solution of the system  $(I + \lambda L_p)f = Y$  based on contour integrals and Krylov subspace methods. The method exploits the sparsity of the Laplacians of each layer and is matrix-free, in the sense that it requires only to compute the matrix-vector product  $L_{\text{sym}}^{(i)} \times \text{vector}$ , without requiring to store the matrices. Thus, when the layers are sparse, the method scales to large datasets. Observe that this is a critical requirement as  $L_p$  is in general a dense matrix, even for very sparse layers, and thus computing and storing  $L_p$  is very prohibitive for large multilayer graphs. We present a method for negative integer values  $p < 0$ , leaving aside the limit case  $p \rightarrow 0$  as it requires a particular treatment. The following is a brief overview of the proposed approach. Further details are available in the supplementary material.

Let  $A_1, \dots, A_T$  be symmetric positive definite matrices,  $\varphi : \mathbb{C} \rightarrow \mathbb{C}$  defined by  $\varphi(z) = z^{1/p}$  and  $L_p = T^{-1/p} \varphi(S_p)$ , where  $S_p = A_1^p + \dots + A_T^p$ . The proposed method consists of three main steps:

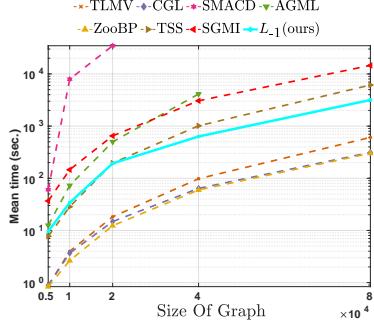


Figure 5: Mean execution time of 10 runs for different methods.  $L_{-1}(\text{ours})$  stands for the power mean Laplacian regularizer together with our proposed matrix-free contour integral based method. We generate multilayer graphs with two layers, each with two classes of same size with parameters  $p_{\text{in}} = 0.05$  and  $p_{\text{in}} = 0.025$  and graphs of sizes  $[0.5, 1, 2, 4, 8] \times 10^4$ . Observe that our matrix free approach for  $L_{-1}$  (solid blue curve) is competitive to state of the art approaches as TSS[28], outperforming AGML[23], SGMI[13] and SMACD[9]. For TLMV[33] and SGMI we use our own implementation.

1. We solve the system  $(I + \lambda L_p)^{-1}Y$  via a Krylov method (e.g. PCG or GMRES) with convergence rate  $O((\frac{\kappa^2 - 1}{\kappa^2})^{h/2})$  [25], where  $\kappa = \lambda_{\max}(L_p)/\lambda_{\min}(L_p)$ . At iteration  $h$ , this method projects the problem onto the Krylov subspace spanned by  $\{Y, \lambda L_p Y, (\lambda L_p)^2 Y, \dots, (\lambda L_p)^h Y\}$ , and efficiently solve the projected problem.
2. The previous step requires the matrix-vector product  $L_p Y = T^{-1/p} \varphi(S_p) Y$  which we compute by approximating the Cauchy integral form of the function  $\varphi$  with the trapezoidal rule in the complex plane [10]. Taking  $N$  suitable contour points and coefficients  $\beta_0, \dots, \beta_N$ , we have

$$\varphi_N(S_p)Y = \beta_0 S_p \operatorname{Im} \left\{ \sum_{i=1}^N \beta_i (z_i^2 I - S_p)^{-1} Y \right\}, \quad (7)$$

which has geometric convergence [10]:  $\|\varphi(S_p)Y - \varphi_N(S_p)Y\| = O(e^{-2\pi^2 N / (\ln(M/m) + 6)})$ , where  $m, M$  are such that  $M \geq \lambda_{\max}(S_p)$  and  $m \leq \lambda_{\min}(S_p)$ .

3. The previous step requires to solve linear systems of the form  $(zI - S_p)^{-1}Y$ . We solve each of these systems via a Krylov subspace method, projecting, at each iteration  $h$ , onto the subspace spanned by  $\{Y, S_p Y, S_p^2 Y, \dots, S_p^h Y\}$ . Since  $S_p = \sum_{i=1}^T A_i^{-|p|}$  this problem reduces to computing  $|p|$  linear systems with  $A_i$  as coefficient matrix, for  $i = 1 \dots, T$ . Provided that  $A_1, \dots, A_T$  are sparse matrices, this is done efficiently using pcg with incomplete Cholesky preconditioners.

Notice that the method allows a high level of parallelism. In fact, the  $N$  (resp.  $p$ ) linear systems solvers at step 2 (resp. 3) are independent and can be run in parallel. Moreover, note that the main task of the method is solving linear systems with Laplacian matrices, which can be solved linearly in the number of edges in the corresponding adjacency matrix. Hence, the proposed approach scales to large sparse graphs and is highly parallelizable. A time execution analysis is provided in Fig 5, where we can see that the time execution of our approach is competitive to the state of the art as TSS[28], outperforming AGML[23], SGMI[13] and SMACD[9].

## 6 Experiments on Real Datasets

In this section we compare the performance of the proposed approach with state of the art methods on real world datasets. We consider the following datasets: *3-sources* [16], which consists of news articles that were covered by news sources BBC, Reuters and Guardian; *BBC*[7] and *BBC Sports*[8] news articles, a dataset of Wikipedia articles with ten different classes [24], the hand written *UCI* digits dataset with six different set of features, and citations datasets *CiteSeer*[17], *Cora*[18] and *WebKB(Texas)*[5]. For each dataset we build the corresponding layer adjacency matrices by taking the symmetric  $k$ -nearest neighbour graph using as similarity measure the Pearson linear correlation, (i.e. we take the  $k$  neighbours with highest correlation), and take the unweighted version of it. Datasets CiteSeer, Cora and WebKB have only two layers, where the first one is a fixed precomputed citation layer, and the second one is the corresponding  $k$ -nearest neighbour graph built from document features.

As **baseline methods** we consider: TSS [28] which identifies an optimal linear combination of graph Laplacians, SGMI [13] which performs label propagation by sparse integration, TLMV [33] which is a weighted arithmetic mean of adjacency matrices, CGL [1] which is a convex combination of the pseudo inverse Laplacian kernel, AGML [23] which is a parameter-free method for optimal graph layer weights, ZooBP [6] which is a fast approximation of Belief Propagation, and SMACD [9] which is a tensor factorization method designed for semi-supervised learning. Finally we set parameters for TSS to ( $c = 10$ ,  $c_0 = 0.4$ ), SMACD ( $\lambda = 0.01$ )<sup>2</sup>, TLMV ( $\lambda = 1$ ), SGMI ( $\lambda_1 = 1$ ,  $\lambda_2 = 10^{-3}$ )

<sup>2</sup>this is the default value in the code released by the authors: <https://github.com/eguirr001/SMACD>

3sources							BBC						
	1%	5%	10%	15%	20%	25%		1%	5%	10%	15%	20%	25%
TLMV	29.8	21.5	<b>20.8</b>	20.3	15.5	16.5	TLMV	<b>29.0</b>	19.3	13.2	11.1	9.3	8.8
CGL	50.2	45.5	36.4	30.6	23.8	19.8	CGL	72.5	52.3	36.1	27.4	22.0	17.1
SMACD	91.5	91.1	91.2	90.9	90.7	91.3	SMACD	74.4	73.5	72.8	72.6	72.5	72.4
AGML	<b>23.9</b>	26.3	33.9	33.3	26.1	22.0	AGML	60.0	34.2	18.6	13.1	11.0	9.5
ZooBP	31.0	21.9	<b>21.3</b>	19.8	15.0	15.3	ZooBP	31.1	20.1	15.0	12.2	10.0	9.1
TSS	29.8	23.9	33.1	34.6	34.8	35.0	TSS	40.4	26.1	20.9	20.1	19.8	19.7
SGMI	34.4	26.6	<b>25.4</b>	24.4	19.1	17.9	SGMI	37.6	28.9	24.9	22.8	20.7	19.3
$L_1$	33.5	23.9	23.4	20.1	15.6	<b>14.6</b>	$L_1$	31.3	22.8	17.4	13.5	10.2	8.9
$L_{-1}$	<b>28.4</b>	<b>20.0</b>	21.8	22.0	17.2	17.9	$L_{-1}$	<b>31.0</b>	<b>17.0</b>	<b>11.5</b>	<b>10.5</b>	<b>9.2</b>	<b>8.7</b>
$L_{-10}$	40.9	29.1	21.9	<b>19.3</b>	<b>14.8</b>	14.7	$L_{-10}$	51.6	26.9	16.6	12.8	10.3	9.5
BBCS							Wikipedia						
	1%	5%	10%	15%	20%	25%		1%	5%	10%	15%	20%	25%
TLMV	25.6	12.6	10.5	7.5	6.4	5.4	TLMV	<b>65.7</b>	<b>56.8</b>	46.4	43.1	40.8	39.2
CGL	79.2	51.6	34.9	23.4	16.5	12.7	CGL	87.3	83.0	82.5	82.2	83.0	83.0
SMACD	77.8	80.6	<b>82.4</b>	96.4	98.4	98.3	SMACD	85.4	85.6	85.4	85.3	86.8	90.0
AGML	34.6	17.4	12.1	<b>7.0</b>	<b>6.0</b>	<b>5.4</b>	AGML	71.3	66.6	48.1	42.1	38.4	37.3
ZooBP	33.8	13.9	11.3	8.8	7.6	6.2	ZooBP	67.6	58.0	47.0	43.8	41.2	39.8
TSS	<b>23.9</b>	13.2	14.1	12.3	13.1	12.2	TSS	87.7	84.7	83.3	81.9	82.3	81.4
SGMI	31.9	19.6	16.6	15.5	14.8	12.1	SGMI	69.3	84.8	84.5	83.8	83.2	82.8
$L_1$	29.9	15.0	13.5	10.6	8.7	7.2	$L_1$	68.2	61.1	53.6	48.3	44.1	42.3
$L_{-1}$	<b>23.8</b>	<b>11.6</b>	<b>8.7</b>	<b>6.3</b>	<b>5.8</b>	<b>5.1</b>	$L_{-1}$	<b>59.1</b>	<b>52.3</b>	<b>40.2</b>	<b>36.3</b>	<b>35.1</b>	<b>34.1</b>
$L_{-10}$	48.7	22.5	14.2	9.1	7.8	6.1	$L_{-10}$	66.9	57.2	43.2	38.7	36.3	34.9
UCI							Citeseer						
	1%	5%	10%	15%	20%	25%		1%	5%	10%	15%	20%	25%
TLMV	28.9	20.4	16.3	14.4	13.7	12.7	TLMV	<b>51.5</b>	39.4	36.5	33.7	31.6	30.3
CGL	81.8	64.0	54.6	49.1	46.7	46.7	CGL	89.3	71.8	58.0	49.8	44.5	40.9
SMACD	73.6	81.0	90.0	90.0	86.2	81.9	SMACD	90.7	90.4	67.0	65.5	66.8	68.9
AGML	<b>25.3</b>	17.2	15.2	13.2	12.5	12.0	AGML	<b>47.3</b>	<b>32.3</b>	<b>29.6</b>	<b>28.2</b>	<b>27.5</b>	<b>27.0</b>
ZooBP	30.8	21.7	17.6	15.1	14.1	13.0	ZooBP	63.6	41.9	38.7	35.8	33.8	32.2
TSS	<b>24.0</b>	17.6	16.6	15.9	15.8	15.6	TSS	58.5	49.5	45.9	42.1	39.8	38.4
SGMI	36.0	44.4	50.9	50.4	50.2	48.8	SGMI	59.4	46.8	44.0	42.3	40.5	39.2
$L_1$	31.3	23.8	18.7	15.6	14.4	13.2	$L_1$	56.3	44.1	41.2	38.5	36.1	34.7
$L_{-1}$	<b>30.5</b>	<b>17.1</b>	<b>13.8</b>	<b>12.6</b>	<b>12.3</b>	<b>11.9</b>	$L_{-1}$	52.4	39.0	35.6	32.6	30.9	29.5
$L_{-10}$	57.0	33.8	23.7	17.6	15.3	13.4	$L_{-10}$	68.6	54.6	48.5	43.0	39.7	37.2
Cora							WebKB						
	1%	5%	10%	15%	20%	25%		1%	5%	10%	15%	20%	25%
TLMV	46.0	34.1	28.8	25.8	22.5	20.6	TLMV	58.6	49.4	45.6	47.2	47.6	48.2
CGL	85.5	70.1	56.5	49.1	44.2	40.0	CGL	80.4	82.4	84.4	86.9	82.7	89.2
SMACD	75.6	76.7	78.7	78.7	81.0	87.1	SMACD	87.3	87.2	87.2	87.4	87.8	87.8
AGML	54.7	36.0	25.4	<b>20.7</b>	<b>18.1</b>	<b>16.5</b>	AGML	56.5	50.3	46.8	44.7	47.6	46.8
ZooBP	54.7	38.0	32.9	30.2	27.6	26.2	ZooBP	52.0	45.0	<b>38.7</b>	38.5	<b>36.4</b>	<b>33.5</b>
TSS	<b>38.8</b>	<b>27.7</b>	<b>24.1</b>	21.5	20.0	19.1	TSS	60.9	51.0	50.5	47.3	49.2	48.7
SGMI	57.3	47.7	43.0	41.8	40.1	38.5	SGMI	<b>44.9</b>	<b>39.7</b>	41.9	<b>34.9</b>	40.3	52.5
$L_1$	50.7	38.2	33.4	31.2	28.2	25.6	$L_1$	58.5	49.0	44.8	44.3	44.5	44.4
$L_{-1}$	<b>43.2</b>	31.8	<b>24.5</b>	21.1	18.8	17.2	$L_{-1}$	<b>49.9</b>	45.5	40.7	39.5	39.9	40.3
$L_{-10}$	62.0	46.3	35.4	29.4	25.2	22.3	$L_{-10}$	52.3	41.9	<b>38.0</b>	38.1	36.8	39.5

Table 2: Experiments in real datasets. Notation: **best** performances are marked with bold fonts and gray background and second best performances with only gray background.

and  $\lambda = 0.1$  for  $L_1$  and  $\lambda = 10$  for  $L_{-1}$  and  $L_{-10}$ . We do not perform cross validation in our experimental setting due to the large execution time in some of the methods here considered. Hence we fix the parameters for each method in all experiments.

We fix nearest neighbourhood size to  $k = 10$  and generate 10 samples of labeled nodes, where the percentage of labeled nodes per class is in the range  $\{1\%, 5\%, 10\%, 15\%, 20\%, 25\%\}$ . The average test errors are presented in table 2, where the **best** (resp. second best) performances are marked with bold fonts and gray background (resp. with only gray background). We can see that the first and second best positions are in general taken by the power mean Laplacian regularizers  $L_1$ ,  $L_{-1}$ ,  $L_{-10}$ , being clear for all datasets except with 3-sources. Moreover we can see that in 77% of all cases  $L_{-1}$  presents either the best or the second best performance, further verifying that our proposed approach based on the power mean Laplacian for semi-supervised learning in multilayer graph is a competitive alternative to state of the art methods<sup>3</sup>.

<sup>3</sup>Communications with the authors of [9] could not clarify the bad performance of SMACD.

**Acknowledgement** P.M and M.H are supported by the DFG Cluster of Excellence “Machine Learning – New Perspectives for Science”, EXC 2064/1, project number 390727645

## References

- [1] A. Argyriou, M. Herbster, and M. Pontil. Combining graph Laplacians for semi-supervised learning. In *NeurIPS*, 2006.
- [2] M. Belkin, I. Matveeva, and P. Niyogi. Regularization and semi-supervised learning on large graphs. In *COLT*, 2004.
- [3] K. V. Bhagwat and R. Subramanian. Inequalities between means of positive operators. *Mathematical Proceedings of the Cambridge Philosophical Society*, 83(3):393–401, 1978.
- [4] O. Chapelle, B. Schölkopf, and A. Zien. *Semi-Supervised Learning*. The MIT Press, 2010.
- [5] M. Craven, D. DiPasquo, D. Freitag, A. McCallum, T. Mitchell, K. Nigam, and S. Slattery. Learning to extract symbolic knowledge from the world wide web. In *AAAI*, 2011.
- [6] D. Eswaran, S. Günnemann, C. Faloutsos, D. Makhija, and M. Kumar. Zoobp: Belief propagation for heterogeneous networks. In *VLDB*, 2017.
- [7] D. Greene and P. Cunningham. Producing accurate interpretable clusters from high-dimensional data. In *PKDD*, 2005.
- [8] D. Greene and P. Cunningham. A matrix factorization approach for integrating multiple data views. In *ECML PKDD*, 2009.
- [9] E. Gujral and E. E. Papalexakis. SMACD: Semi-supervised multi-aspect community detection. In *SDM*, 2018.
- [10] N. Hale, N. J. Higham, and L. N. Trefethen. Computing  $A^\alpha$ ,  $\log(A)$ , and related matrix functions by contour integrals. *SIAM Journal on Numerical Analysis*, 46(5):2505–2523, 2008.
- [11] S. Heimlicher, M. Lelarge, and L. Massoulié. Community detection in the labelled stochastic block model. *arXiv:1209.2910*, 2012.
- [12] V. Kanade, E. Mossel, and T. Schramm. Global and local information in clustering labeled block models. *IEEE Transactions on Information Theory*, 62(10):5906–5917, 2016.
- [13] M. Karasuyama and H. Mamitsuka. Multiple graph label propagation by sparse integration. *IEEE Transactions on Neural Networks and Learning Systems*, 24(12):1999–2012, 2013.
- [14] T. Kato, H. Kashima, and M. Sugiyama. Robust label propagation on multiple networks. *Transactions on Neural Networks*, 20(1):35–44, Jan. 2009.
- [15] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- [16] J. Liu, C. Wang, J. Gao, and J. Han. Multi-view clustering via joint nonnegative matrix factorization. In *SDM*, 2013.
- [17] Q. Lu and L. Getoor. Link-based classification. In *ICML*, 2003.
- [18] A. K. McCallum, K. Nigam, J. Rennie, and K. Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3(2):127–163, 2000.
- [19] P. Mercado, A. Gautier, F. Tudisco, and M. Hein. The power mean Laplacian for multilayer graph clustering. In *AISTATS*, 2018.
- [20] P. Mercado, F. Tudisco, and M. Hein. Clustering signed networks with the geometric mean of Laplacians. In *NeurIPS*. 2016.
- [21] P. Mercado, F. Tudisco, and M. Hein. Spectral clustering of signed graphs via matrix power means. In *ICML*, 2019.

- [22] E. Mossel and J. Xu. Local algorithms for block models with side information. In *ITCS*, 2016.
- [23] F. Nie, J. Li, and X. Li. Parameter-free auto-weighted multiple graph learning: A framework for multiview clustering and semi-supervised classification. In *IJCAI*, 2016.
- [24] N. Rasiwasia, J. Costa Pereira, E. Coviello, G. Doyle, G. R. Lanckriet, R. Levy, and N. Vasconcelos. A new approach to cross-modal multimedia retrieval. In *ACM Multimedia*, 2010.
- [25] Y. Saad and M. H. Schultz. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal on Scientific and Statistical Computing*, 7(3):856–869, 1986.
- [26] A. Saade, F. Krzakala, M. Lelarge, and L. Zdeborová. Fast randomized semi-supervised clustering. *Journal of Physics: Conference Series*, 1036:012015, 2018.
- [27] A. Subramanya and P. P. Talukdar. *Graph-Based Semi-Supervised Learning*. Morgan & Claypool Publishers, 2014.
- [28] K. Tsuda, H. Shin, and B. Schölkopf. Fast protein classification with multiple networks. *Bioinformatics*, 21(2):59–65, 2005.
- [29] K. Viswanathan, S. Sachdeva, A. Tomkins, and S. Ravi. Improved semi-supervised learning with multiple graphs. In *AISTATS*, 2019.
- [30] Z. Yang, W. W. Cohen, and R. Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *ICML*, 2016.
- [31] J. Ye and L. Akoglu. Robust semi-supervised learning on multiple networks with noise. In *PKDD*, 2018.
- [32] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf. Learning with local and global consistency. In *NeurIPS*, 2003.
- [33] D. Zhou and C. J. Burges. Spectral clustering and transductive learning with multiple views. In *ICML*, 2007.
- [34] X. Zhu and Z. Ghahramani. Learning from labeled and unlabeled data with label propagation. Technical report, 2002.
- [35] X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*, 2003.