

Influence of Graph Construction on Semi-supervised Learning

Celso André R. de Sousa, Solange O. Rezende,
and Gustavo E.A.P.A. Batista

Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo,
Campus de São Carlos, Brazil
{sousa,solange,gbatista}@icmc.usp.br

Abstract. A variety of graph-based semi-supervised learning (SSL) algorithms and graph construction methods have been proposed in the last few years. Despite their apparent empirical success, the field of SSL lacks a detailed study that empirically evaluates the influence of graph construction on SSL. In this paper we provide such an experimental study. We combine a variety of graph construction methods as well as a variety of graph-based SSL algorithms and empirically compare them on a number of benchmark data sets widely used in the SSL literature. The empirical evaluation proposed in this paper is subdivided into four parts: (1) best case analysis; (2) classifiers' stability evaluation; (3) influence of graph construction; and (4) influence of regularization parameters. The purpose of our experiments is to evaluate the trade-off between classification performance and stability of the SSL algorithms on a variety of graph construction methods and parameter values. The obtained results show that the mutual k -nearest neighbors (mutKNN) graph may be the best choice for adjacency graph construction while the RBF kernel may be the best choice for weighted matrix generation. In addition, mutKNN tends to generate smoother error surfaces than other adjacency graph construction methods. However, mutKNN is unstable for a relatively small value of k . Our results indicate that the classification performance of the graph-based SSL algorithms are heavily influenced by the parameters setting and we found no evident explorable pattern to relay to future practitioners. We discuss the consequences of such instability in research and practice.

Keywords: Semi-supervised learning, graph-based methods, experimental study, classification.

1 Introduction

Semi-supervised learning (SSL) has gained increased attention in the last few years [3,15]. Among all SSL algorithms, graph-based methods are widely used because the weighted graph may approximate the low dimensional manifold in which the data should lie. The research community has proposed a variety of graph-based SSL algorithms [1,8,14,16] as well as a variety of graph construction

methods [5,8,13]. Despite its increasing popularity, the SSL literature lacks a comprehensive and unbiased empirical study that shows the influence that graph construction methods have in both classification performance and stability of the graph-based SSL algorithms.

1.1 Contributions

In this paper, we provide a detailed empirical comparison of the state-of-the-art, graph-based SSL algorithms combined with a variety of graph construction methods. The empirical analysis proposed in this paper is subdivided into four parts as follows:

Best case analysis. We evaluate the best error rates of each combination of SSL algorithm and graph construction method for a number of sparsification parameter values. Although this is the most common approach to evaluate SSL algorithms in the literature [3], this empirical setting alone may not provide all the necessary information to choose the best classifiers for real applications. For instance, stable classifiers may be preferable over classifiers which are able to provide excellent performance for a very narrow range of parameter values and mediocre performance for the remaining values;

Classifiers' stability evaluation. We evaluate the stability of the SSL algorithms combined with the graph construction methods as we vary the value of the sparsification parameter. As we mentioned before, this analysis is important because a classifier may achieve the best overall classification performance for a very narrow range of the parameter values. Then, this analysis is an invaluable tool to identify which classifiers provide a good trade-off between classification performance and stability;

Influence of graph construction. We also evaluate the graph construction methods combined with the SSL algorithms over a wide range of sparsification parameter values. We want to verify: (1) how the graph construction methods affect the classification performance of each SSL algorithm and (2) the stability of the graph construction methods as we vary the sparsification parameter values. For the classifiers that have at least one regularization parameter, we fixed the regularization parameter(s) with the value that achieved the best average error rate and then varied the sparsification parameter value;

Influence of regularization parameters. We evaluate the error surfaces generated by the SSL algorithms that have regularization parameters. We first chose the sparsification parameter that achieved the best average error rate and then we varied the regularization parameters of the SSL algorithms.

The obtained results show that the mutual k -nearest neighbors (mutKNN) graph may be the best choice for adjacency graph construction while the RBF kernel may be the best choice for weighted matrix generation. In addition, mutKNN tends to generate smoother error surfaces than other adjacency graph construction methods. However, mutKNN is unstable for a relatively small value of k .

Our results indicate that the classification performance of the graph-based SSL algorithms are heavily influenced by *internal* parameters (such as regularization parameters) and *external* parameters (such as the number of neighbors in a k -nearest neighbor graph). Such variability showed no evident explorable pattern to relay to future practitioners. In addition, the SSL assumption that only a very restricted set of labeled examples exists may make parameter estimation techniques commonly used in classification unfeasible.

We believe that our results have two major consequences:

For practitioners. Given a data set, it is difficult to recommend an SSL algorithm, a graph sparsification parameter value or a regularization parameter value that is expected to provide good classification performance. As the number of labeled examples is usually very restricted in SSL applications, the practitioner has no tools to make an informed choice of these parameter values. As we will show, an incorrect choice of the parameter values may seriously affect the classification results;

For researchers. Changes in the parameter values also cause changes in the relative ranking among the classifiers. It means that for a specific data set several methods may figure as the best classifier for a certain range of parameter values. This is a serious issue since the empirical evidence that one method outperforms the competitors might be confirmed only for a restricted set of the parameter values. In addition, this performance variability may hinder the reproduction of the experimental results for papers that do not clearly report every parameter value used in the empirical evaluation.

1.2 Outline

The remainder of this paper is organized as follows. Section 2 describes the notation used throughout the paper and revises the graph construction methods. Section 3 revises the state-of-the-art, graph-based SSL algorithms. Section 4 empirically evaluates the graph construction methods combined with the graph-based SSL algorithms. Finally, Section 5 concludes the paper and suggests directions for future research.

2 Graph Construction

In this section we revise widely used methods to generate sparse weighted graphs, which are frequently considered the heart of graph-based SSL [15]. Section 2.1 describes the notation used throughout the paper. Section 2.2 revises approaches used to generate a sparse undirected¹ graph (or adjacency matrix) from the training sample. Section 2.3 revises approaches used to generate a weighted matrix from the sparse graph.

¹ This paper focus on undirected graphs, which are commonly used in SSL [15].

2.1 Notation and Preliminaries

Consider a training sample $\mathcal{X} := \{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^d$ in which the first l examples are labeled, i.e., \mathbf{x}_i has label $y_i \in \mathbb{N}_c$ where $\mathbb{N}_p := \{i \in \mathbb{N}^* | 1 \leq i \leq p\}$ with $p \in \mathbb{N}^*$ and c being the number of classes. Let $u := n - l$ be the amount of unlabeled examples and $\mathbf{Y} \in \mathbb{B}^{n \times c}$ be a label matrix in which $\mathbf{Y}_{ij} = 1$ if and only if \mathbf{x}_i has label $y_i = j$. Consider an undirected graph $\mathcal{G} := (\mathcal{X}, \mathcal{E})$ in which each \mathbf{x}_i is a node of \mathcal{G} . Let $\mathcal{N}_i \subset \mathcal{X}$ be the set of neighbors of \mathbf{x}_i and \mathbf{x}_{i_k} the k -th nearest neighbor of \mathbf{x}_i . In order to generate a sparse weighted matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ from \mathcal{G} one uses a similarity function $\mathcal{K} : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$ to compute the weights \mathbf{W}_{ij} .

The graph Laplacians are important tools for machine learning. The *combinatorial* Laplacian is defined by $\mathbf{\Delta} := \mathbf{D} - \mathbf{W}$ where $\mathbf{D} := \text{diag}(\mathbf{W}\mathbf{1}_n)$ such that $\mathbf{1}_n$ is an n -dimensional 1-entry vector. The *normalized* Laplacian is defined by $\mathbf{L} := \mathbf{I}_n - \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$ where \mathbf{I}_n is the n -by- n identity matrix.

All matrices can be subdivided into labeled and unlabeled submatrices. Let $\mathbf{F} \in \mathbb{R}^{n \times c}$ be the output of a given graph-based SSL algorithm. The \mathbf{F} and \mathbf{Y} matrices are subdivided into two submatrices while all others are subdivided into four submatrices. For instance:

$$\mathbf{W} := \begin{bmatrix} \mathbf{W}_{\mathcal{L}\mathcal{L}} & \mathbf{W}_{\mathcal{L}\mathcal{U}} \\ \mathbf{W}_{\mathcal{U}\mathcal{L}} & \mathbf{W}_{\mathcal{U}\mathcal{U}} \end{bmatrix} \quad \mathbf{Y} := \begin{bmatrix} \mathbf{Y}_{\mathcal{L}} \\ \mathbf{Y}_{\mathcal{U}} \end{bmatrix}$$

where $\mathbf{W}_{\mathcal{L}\mathcal{L}} \in \mathbb{R}^{l \times l}$ and $\mathbf{Y}_{\mathcal{L}} \in \mathbb{B}^{l \times c}$ are the submatrices of \mathbf{W} and \mathbf{Y} , respectively, on labeled examples, and so on. By definition, $\mathbf{Y}_{\mathcal{U}}$ is an $u \times c$ null matrix. This paper focus on the multi-class problem; hence, $\mathbf{Y}_{\mathcal{L}}\mathbf{1}_c = \mathbf{1}_l$.

2.2 Adjacency Graph Construction

The adjacency graph construction process generates a graph \mathcal{G} (or adjacency matrix \mathbf{A}) from \mathcal{X} using a distance function $\Psi : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$. Let $\mathbf{\Psi} \in \mathbb{R}^{n \times n}$ be a distance matrix in which $\mathbf{\Psi}_{ij} := \Psi(\mathbf{x}_i, \mathbf{x}_j)$ and $\mathbf{A} \in \mathbb{B}^{n \times n}$ be an adjacency matrix² in which $\mathbf{A}_{ij} = 1$ if and only if $\mathbf{x}_j \in \mathcal{N}_i$. We now describe the two most used adjacency graph construction methods for graph-based learning.

ϵ -neighborhood ($\epsilon\mathbf{N}$). There exists an undirected edge between \mathbf{x}_i and \mathbf{x}_j in an $\epsilon\mathbf{N}$ graph if and only if $\Psi(\mathbf{x}_i, \mathbf{x}_j) \leq \epsilon$ where $\epsilon \in \mathbb{R}_+^*$ is a free parameter. In general, $\epsilon\mathbf{N}$ graphs are not widely used in practical situations because they can generate graphs with many disconnected components for an improper value of ϵ . Due to this fact, we did not use the $\epsilon\mathbf{N}$ graph in our experiments.

k -nearest neighbors ($k\mathbf{NN}$). There exists an edge from \mathbf{x}_i to \mathbf{x}_j if and only if \mathbf{x}_j is one of the k closest examples of \mathbf{x}_i . Because the adjacency matrix of a $k\mathbf{NN}$ graph may not be symmetric, three strategies are commonly used to symmetrize it: *mutual $k\mathbf{NN}$* (mutKNN), which generates $\widehat{\mathbf{A}} = \min(\mathbf{A}, \mathbf{A}^\top)$; *symmetric $k\mathbf{NN}$* (symKNN), which generates $\widehat{\mathbf{A}} = \max(\mathbf{A}, \mathbf{A}^\top)$; and *symmetry-favored $k\mathbf{NN}$* (symFKNN) [8], which generates $\widehat{\mathbf{A}} = \mathbf{A} + \mathbf{A}^\top$ (a non-binary adjacency matrix).

² Non-binary adjacency matrices may also be applied.

2.3 Weighted Matrix Generation

Given an adjacency matrix \mathbf{A} , we generate a sparse weighted matrix \mathbf{W} using a similarity function $\mathcal{K} : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$. We describe three widely used approaches to generate \mathbf{W} . Two of them, RBF kernel and similarity function of Hein & Maier [5], define the \mathbf{W} matrix using the relation $\mathbf{W}_{ij} = \mathbf{A}_{ij}\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$. The third approach, based on local reconstruction minimization [13], generates a sparse weighted matrix \mathbf{W} , not necessarily symmetric, without an explicit \mathcal{K} .

RBF kernel. The RBF (or Gaussian) kernel computes the similarity between \mathbf{x}_i and \mathbf{x}_j by $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) := \exp(-\Psi^2(\mathbf{x}_i, \mathbf{x}_j)/(2\sigma^2))$ in which $\sigma \in \mathbb{R}_+^*$ is the kernel bandwidth parameter.

Similarity function of Hein & Maier [5] (HM). Given a function $\psi(\cdot, \cdot)$ in which $\psi(\mathbf{x}_i, k) := \Psi(\mathbf{x}_i, \mathbf{x}_{i_k})$ with $k \in \mathbb{N}^*$, the HM similarity function is defined by $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) := \exp\left(-\Psi^2(\mathbf{x}_i, \mathbf{x}_j)/(\max\{\psi(\mathbf{x}_i, k), \psi(\mathbf{x}_j, k)\})^2\right)$. This is an RBF kernel with an adaptive kernel size.

Local Linear Embedding (LLE). The LLE approach [13] generates the \mathbf{W} matrix by solving the following optimization problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{n \times n}} \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{\mathbf{x}_j \in \mathcal{N}_i} \mathbf{W}_{ij} \mathbf{x}_j \right\|_2^2 \quad \text{s.t.} \quad \mathbf{W} \mathbf{1}_n = \mathbf{1}_n, \quad \mathbf{W} \geq 0 \quad (1)$$

The symbol $\|\cdot\|_2$ represents the l_2 -norm.

3 Label Diffusion

Given a weighted matrix \mathbf{W} , a graph-based SSL algorithm uses \mathbf{W} and the label matrix \mathbf{Y} to generate the output matrix \mathbf{F} by label diffusion in the weighted graph. We now revise the state-of-the-art graph-based SSL algorithms used in our empirical comparison. We should note that these algorithms have an intrinsic condition to classify all unlabeled examples in \mathcal{X} , which frequently is not explicit in the literature. Assumption 1 describes this condition.

Assumption 1. *Each unlabeled example is on a connected subgraph in which there exists at least one labeled example.*

Gaussian Random Fields (GRF). The GRF algorithm [16] solves the optimization problem $\mathbf{F} = \arg \min_{\mathbf{F} \in \mathbb{R}^{n \times c}} \text{tr}(\mathbf{F}^\top \Delta \mathbf{F})$ s.t. $\mathbf{F}_\mathcal{L} = \mathbf{Y}_\mathcal{L}$, which gives the closed-form solution $\mathbf{F}_\mathcal{U} = -\Delta_{\mathcal{U}\mathcal{U}}^{-1} \Delta_{\mathcal{U}\mathcal{L}} \mathbf{Y}_\mathcal{L}$.

Local and Global Consistency (LGC). The LGC algorithm [14] solves the optimization problem $\mathbf{F} = \arg \min_{\mathbf{F} \in \mathbb{R}^{n \times c}} \text{tr}(\mathbf{F}^\top \mathbf{L} \mathbf{F} + \mu(\mathbf{F} - \mathbf{Y})^\top (\mathbf{F} - \mathbf{Y}))$, which gives the closed-form solution $\mathbf{F} = (\mathbf{I}_n + \mathbf{L}/\mu)^{-1} \mathbf{Y}$.

Laplacian Regularized Least Squares (LapRLS). The LapRLS algorithm [1] minimizes the following regularization framework:

$$\min_{f \in \mathcal{H}_\mathcal{K}} \frac{1}{l} \sum_{i=1}^l \mathcal{V}(\mathbf{x}_i, y_i, f) + \gamma_A \|f\|_{\mathcal{H}_\mathcal{K}} + \gamma_I \mathbf{f}^\top \Delta \mathbf{f} \quad (2)$$

where $\mathcal{V}(\mathbf{x}_i, y_i, f) = (y_i - f(\mathbf{x}_i))^2$, $\mathcal{H}_{\mathcal{K}}$ is the *Reproducing Kernel Hilbert Space (RKHS)* for the kernel \mathcal{K} , $\mathbf{f} := [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^\top \in \mathbb{R}^n$, $\|\cdot\|_{\mathcal{H}_{\mathcal{K}}}$ is the norm in $\mathcal{H}_{\mathcal{K}}$, and γ_A and γ_I are the regularization parameters. Let $\mathbf{y} := [y_1, \dots, y_l, 0, \dots, 0] \in \mathbb{R}^n$ be the label vector in which $y_i \in \{-1, +1\}$ and $\mathbf{K} \in \mathbb{R}^{n \times n}$ a gram matrix such that $\mathbf{K}_{ij} := \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$. Due to the *Representer Theorem* in [1], the solution of (2) can be written as an expansion of kernel functions over both labeled and unlabeled examples, i.e., $f(\mathbf{x}) = \sum_{i=1}^n \mathcal{K}(\mathbf{x}, \mathbf{x}_i) \alpha_i$ with $\alpha \in \mathbb{R}^n$. Solving (2) using this expansion, we get $\alpha = (\mathbf{JK} + \gamma_A \mathbf{I}_n + \gamma_I l \Delta \mathbf{K})^{-1} \mathbf{y}$ where $\mathbf{J} := \text{diag}([1, \dots, 1, 0, \dots, 0]^\top)$ whose first l diagonal entries are 1 and the rest 0.

Laplacian Support Vector Machine (LapSVM). The LapSVM algorithm [1] minimizes the problem in (2) with $\mathcal{V}(\mathbf{x}_i, y_i, f) = \max(0, 1 - y_i f(\mathbf{x}_i))$. Solving (2) using the expansion $f(\mathbf{x}) = \sum_{i=1}^n \mathcal{K}(\mathbf{x}, \mathbf{x}_i) \alpha_i$, we get the solution $\alpha = \frac{1}{2} (\gamma_A \mathbf{I}_n + \gamma_I \Delta \mathbf{K})^{-1} \bar{\mathbf{J}}^\top \bar{\mathbf{Y}} \beta^*$ where $\bar{\mathbf{J}} := [\mathbf{I}_l \ \mathbf{O}_{l \times u}]$ such that $\mathbf{O}_{l \times u}$ is an $l \times u$ null matrix, $\bar{\mathbf{Y}} := \text{diag}([y_1, \dots, y_l]^\top)$, and $\beta^* \in \mathbb{R}^l$ is given by

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^l} \mathbf{1}_l^\top \beta - \frac{1}{2} \beta^\top \mathbf{Q} \beta \quad \text{s.t.} \quad \mathbf{y}^\top \beta = 0, \quad 0 \leq \beta \leq \frac{1}{l}$$

such that $\mathbf{Q} = \frac{1}{2} \bar{\mathbf{Y}} \bar{\mathbf{J}} \mathbf{K} (\gamma_A \mathbf{I}_n + \gamma_I \Delta \mathbf{K})^{-1} \bar{\mathbf{J}}^\top \bar{\mathbf{Y}}$.

Robust Multi-class Graph Transduction (RMGT). The RMGT algorithm [8] solves the convex optimization problem $\mathbf{F} = \arg \min_{\mathbf{F} \in \mathbb{R}^{n \times c}} \text{tr}(\mathbf{F}^\top \Delta \mathbf{F})$ s.t. $\mathbf{F}_{\mathcal{L}} = \mathbf{Y}_{\mathcal{L}}$, $\mathbf{F} \mathbf{1}_c = \mathbf{1}_n$, $\mathbf{F}^\top \mathbf{1}_n = n \boldsymbol{\omega}$ where $\boldsymbol{\omega} \in \mathbb{R}^c$ is the class prior probabilities. The solution of this optimization problem is given by:

$$\mathbf{F}_U = -\Delta_{UU}^{-1} \Delta_{UL} \mathbf{Y}_{\mathcal{L}} + \frac{\Delta_{UU}^{-1} \mathbf{1}_U}{\mathbf{1}_U^\top \Delta_{UU}^{-1} \mathbf{1}_U} (n \boldsymbol{\omega}^\top - \mathbf{1}_l^\top \mathbf{Y}_{\mathcal{L}} + \mathbf{1}_U^\top \Delta_{UU}^{-1} \Delta_{UL} \mathbf{Y}_{\mathcal{L}})$$

4 Experimental Evaluation

In this section we provide a detailed empirical comparison of the graph-based SSL algorithms described in Section 3 combined with the graph construction methods described in Section 2 on a number of benchmark data sets. The objective of these experiments is to evaluate the influence that graph construction methods have in the classifiers’ performance. We performed experiments in a transductive setting using different sets of labeled and unlabeled examples in each execution.

For a fair comparison and ease of reproducibility, we used the source code of the authors of the algorithms when possible. As some authors implemented their methods in Matlab, we used the matlabcontrol³ library to link the Matlab code and Java. Due to reasons concerning reproducibility, all source codes and data sets used in our experiments are freely available⁴.

³ <https://code.google.com/p/matlabcontrol/downloads/list>

⁴ <http://www.icmc.usp.br/~gbatista/ECML2013>

4.1 Data Sets

We used in our experiments the USPS, COIL₂, DIGIT-1, G-241C, G-241N, and TEXT data sets. These data sets are freely available⁵ and very popular in the SSL literature [3]. USPS and DIGIT-1 are data sets for digit recognition, TEXT is a data set for text classification, G-241N and G-241C are data sets for classification of Gaussian distributions, and COIL₂ is a data set for image classification. We used the data splits of 10 labeled examples suggested in [3].

We run *principal component analysis* (PCA) to reduce the dimensionality of the data sets. In high-dimensional data, the distance to the nearest neighbor approaches the distance of the farthest neighbor [2]. It degenerates the quality of the graph and possibly decreases the classification performance of the SSL algorithms. After some preliminary experimental evaluation, we decided to reduce the dimensionality of the data to 50 features using the *Matlab Toolbox for Dimensionality Reduction*⁶ library. We did not run PCA only on the TEXT data set to maintain the sparseness property of these data.

4.2 Empirical Setup

In this section, we describe the experimental design decisions that we have taken in our experiments in order to facilitate the reproduction of our results.

Distance functions. Due to its high popularity in the text classification literature, we used the cosine distance in the experiments using the TEXT data set. The cosine distance is defined as $\Psi(\mathbf{x}_i, \mathbf{x}_j) = 1 - \langle \mathbf{x}_i, \mathbf{x}_j \rangle_d / (\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2)$ where $\langle \cdot, \cdot \rangle_d$ is the inner product of vectors in \mathbb{R}^d . For all other data sets we used the l_2 norm as a distance function.

Graph Laplacians. Since the normalized Laplacian \mathbf{L} may lead to better empirical results in comparison with the combinatorial Laplacian $\mathbf{\Delta}$ [7], we used \mathbf{L} instead of $\mathbf{\Delta}$ in the formulation of the graph-based SSL algorithms. We obtained poor results using \mathbf{L} in the RMGT algorithm during preliminary experiments; therefore, we report the results of RMGT using $\mathbf{\Delta}$. In preliminary experiments, we observed some errors using RMGT in the COIL₂ data set. These errors occurred because at least one of the eigenvalues of the graph Laplacian was equal to (or approximately) zero. In an attempt to avoid numerical instabilities while solving linear systems using the graph Laplacians, we generated the combinatorial Laplacian as $\mathbf{\Delta} = \gamma \mathbf{D} - \mathbf{W}$ and the normalized Laplacian as $\mathbf{L} = \gamma \mathbf{I}_n - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$ where a small $\gamma > 1$ is used to increase the eigenvalues of the graph Laplacians. In our experiments, we set $\gamma = 1.01$.

Mutual k NN. The procedure $\hat{\mathbf{A}} = \min(\mathbf{A}, \mathbf{A}^\top)$ may generate a graph with isolated vertices. It may degenerate the output of the SSL algorithms because the label diffusion process could not be effective. In an attempt to avoid this

⁵ <http://olivier.chapelle.cc/ssl-book/benchmarks.html>.

⁶ <http://homepage.tudelft.nl/19j49/>

[Matlab_Toolbox_for_Dimensionality_Reduction.html](http://homepage.tudelft.nl/19j49/Matlab_Toolbox_for_Dimensionality_Reduction.html).

problem, we created an undirected edge between each isolated vertex and its nearest neighbor. Other strategies may also be applied as well [12].

LLE. We used the *Local Anchor Embedding* (LAE) method [9]⁷ to solve the optimization problem in (1). LLE is an example of LAE if we generate a bipartite graph whose “anchor” points are exactly the training examples. Since LLE may not generate a symmetric weighted matrix, we symmetrize the output matrix of LLE, \mathbf{W}_{LLE} , as $\mathbf{W} = \frac{1}{2} (\mathbf{W}_{LLE} + \mathbf{W}_{LLE}^\top)$.

SymFKNN + LLE. Because the adjacency matrix of the symFKNN graph is non-binary, we compute $\widehat{\mathbf{W}} = \mathbf{W}_{LLE} \odot \mathbf{A}$ where \odot is the Hadamard product. Then, we generate $\mathbf{W} = \frac{1}{2} (\widehat{\mathbf{W}} + \widehat{\mathbf{W}}^\top)$.

LapSVM. We run LapSVM using the source code in [11]⁸. We trained LapSVM using Newton’s method, which gave better results than the preconditioned conjugate gradient method during preliminary experiments.

LapRLS. We used the multi-class version of LapRLS; hence, we compute $\boldsymbol{\alpha}$ as $\boldsymbol{\alpha} = (\mathbf{JK} + \gamma_A \mathbf{I}_n + \gamma_I \mathbf{l} \Delta \mathbf{K})^{-1} \mathbf{Y}$ and get the output matrix $\mathbf{F} = \mathbf{K} \boldsymbol{\alpha}$.

Classification. In order to classify the unlabeled examples, we used the *class mass normalization* (CMN) procedure [16]. This is an useful procedure when we are dealing with data sets with imbalanced labels. We obtained poor results using CMN in RMGT; therefore, we report the results for RMGT using the *argmax* operator. We report the results for GRF, LGC, and LapRLS using CMN while the results for LapSVM are reported using the *sign* function. For GRF, we computed CMN using \mathbf{F}_U instead of \mathbf{F} , as suggested in [16].

4.3 Parameter Setting

We now describe the parameter setting used in our experimental evaluation.

SymKNN, mutKNN, and symFKNN. The sparsification parameter k was chosen at the range $\{4, 6, 8, \dots, 40\}$.

RBF kernel. Because it is not straightforward to find an adequate value for the kernel bandwidth σ when labeled examples are scarce, we estimate its value by $\sigma = \sum_{i=1}^n \Psi(\mathbf{x}_i, \mathbf{x}_{i_k}) / (3n)$, as suggested in [6].

Gram matrix. We generated the gram matrix \mathbf{K} using the RBF kernel. We used the same distance function $\Psi(\cdot, \cdot)$, the sparsification parameter k , and the kernel bandwidth σ used during graph construction to compute \mathbf{K} .

LGC. The regularization parameter μ in the LGC framework was chosen at range $\{0.01, 0.05, 0.1, 0.5, 1, 2, 5, 10, 50, 100\}$.

LapRLS and LapSVM. The regularization parameters γ_A and γ_I were chosen at range $\{10^{-6}, 10^{-4}, 10^{-2}, 10^{-1}, 1, 10, 100\}$, as suggested in [11]. All other parameters were set to their default values.

RMGT. For the RMGT algorithm, we assumed a uniform class distribution, i.e., we set $\boldsymbol{\omega} = \mathbf{1}_c / c$ instead of using the class prior probabilities, as suggested in [8]. We achieved better results in preliminary experiments using

⁷ <http://www.ee.columbia.edu/ln/dvmm/downloads/WeiGraphConstructCode/dlform.htm>.

⁸ <http://www.dii.unisi.it/~melacci/lapsvmp/index.html>.

the uniform class distribution in most data sets; therefore, we report the results for RMGT using this setting for all data sets, excluding USPS. For the USPS data set, we used the class prior probabilities, which achieved the best results.

4.4 Analysis of the Results

In this section we analyze the obtained results. Our empirical analysis is subdivided into four parts: (1) best cases analysis; (2) graph-based SSL algorithm comparison; (3) influence of graph construction on SSL; and (4) influence of regularization parameters on the classifiers' performance.

Best case analysis. Table 1 shows the obtained results for the best case analysis. Each numerical result in this table is the lowest average error rate obtained by a combination of an SSL algorithm, a graph construction method and a data set for all parameter values (sparsification and regularization, if applicable), as described in Section 4.3. The four worst results obtained by an SSL algorithm in each data set have a grey background while the best one is in bold. The best overall result for each data set is boxed.

We can see in Table 1 that the symKNN-LLE and symFKNN-LLE graphs may not be adequate for GRF, LGC, and LapRLS because they achieved unsatisfactory results in all data sets. We also see that mutKNN outperformed the symKNN and symFKNN graphs in most situations, independent of the weighted matrix generation method or the SSL algorithm used. Therefore, for the data sets considered in this study, mutKNN presented the best performance among all adjacency graph construction methods.

We ran the Friedman's test⁹ with Nemenyi's post test using a confidence level of 0.05 to statistically compare the performance of the graph construction methods. Table 2 shows the average rankings. The best rankings are marked in bold face and the results that were outperformed by the best ranked method are marked with grey background. We can see that symFKNN-RBF and mutKNN-RBF obtained the best rankings for most SSL algorithms. However, the statistical test found significant differences for only 7 cases.

After analyzing the classifiers, we see that RMGT achieved the best overall classification performance in 4 out of 6 data sets. Although RMGT achieved satisfactory results on most data sets, it did not perform well on the USPS data set.

Classifiers' stability evaluation. As we mentioned earlier, the best case analysis does not allow us to investigate the stability of the classifiers. In this analysis, we investigate the stability of the SSL algorithms as we vary the graph sparsification parameter value. Due to space restrictions and because the mutKNN-RBF graph achieved the best overall classification performance in the best case analysis, we show here only the results obtained with the mutKNN-RBF graph. The interested reader will find the results for other graph construction methods on the paper's website.

⁹ See [4] and references therein for a review on statistical tests for machine learning.

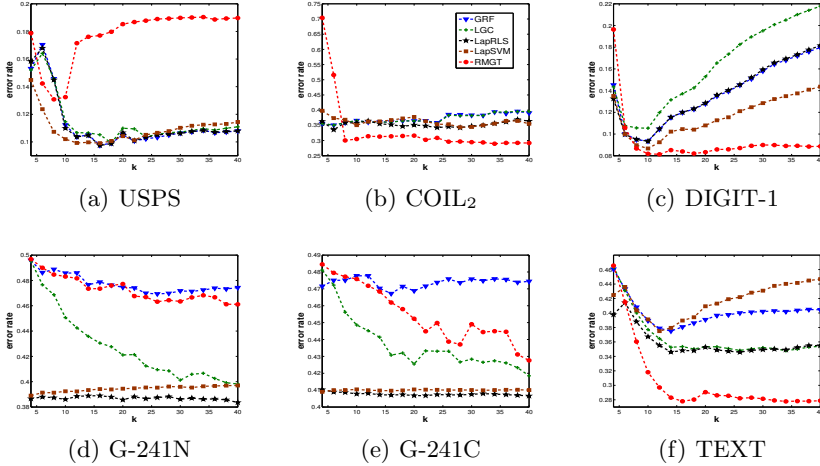


Fig. 1. Average error rates of the SSL algorithms using the mutKNN-RBF graph

Fig. 1 shows the results for this empirical analysis using the mutKNN-RBF graph as we vary the sparsification parameter value. Notice that the legend for all graphics in Fig. 1 can be found in Fig. 1(b). The RMGT algorithm achieved good classification performance and stability on the COIL₂, TEXT, and DIGIT-1 data sets when $k \geq 14$. However, RMGT was generally the worst classifier for the USPS data set and the second worst for the G-241N and G-241C data sets. Moreover, RMGT appears to be unstable for relatively small values of k . For instance, the instability of RMGT is evidenced in the COIL₂ data set for $k \leq 6$ while all other classifiers achieved satisfactory results with this setting.

LapRLS and LapSVM achieved exceptional stability on the G-241C and G-241N data sets. Due to this high stability, we suppose that LapRLS and LapSVM may be the best SSL algorithms for classification of Gaussian distributions. We also note in Fig. 1 that the assumption that sparse graphs give better results than dense graphs may not necessarily be true. For instance, the results for the GRF, LGC, and RMGT algorithms on the G-241C and G-241N data sets using dense graphs are better than those for sparse graphs. In addition, the results for all SSL algorithms on the TEXT data set for relatively small values of k are not satisfactory while the results for the LGC, LapRLS, and RMGT algorithms with dense graphs are.

Influence of graph construction. We now evaluate how different graphs can influence the classification performance of the SSL algorithms. Once again, we perform this analysis as we vary the sparsification parameter value in order to analyze the stability of the graph construction methods combined with the SSL algorithms. Due to lack of space, we only present the plots for the USPS data set in Fig. 2. Once again, we invite the interested reader to check the paper’s website. It is clear from Fig. 2 that the results show

a lot of variability. For a given classifier, we can observe that several graph construction methods figure among the best and the worst method as we vary the value of k . The variability problem is more intense for small values of k , specially $k \leq 14$. This seems to be a permissive problem since small values of k performed better for this specific data set, but too small values might greatly degrade the classifiers' performance.

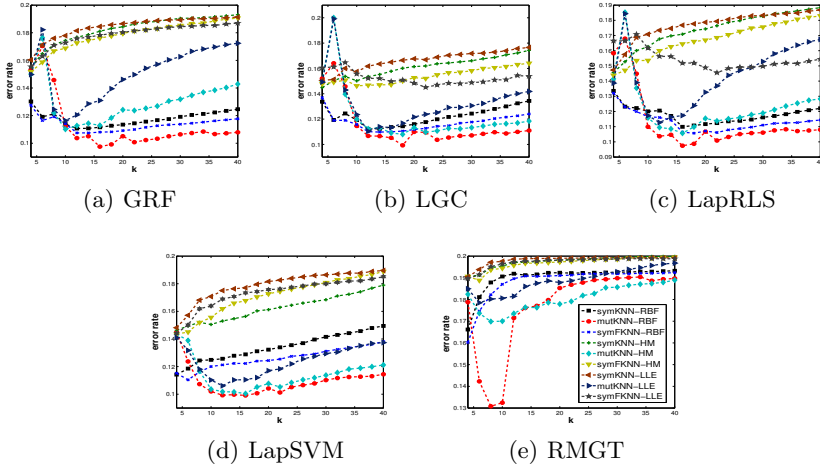


Fig. 2. Average error rates of the graph construction methods on the USPS data set

The USPS data set is an excellent example of the high influence of the graph construction methods and the sparsification parameter values over the SSL algorithms. As we vary the k parameter, the classifiers' performance vary significantly; some of them in a range of almost 10%. Such performance variation is certainly a concern for the practitioner, who would have difficulties in finding a parameter setting that guarantees a good classification performance. Moreover, such high variability causes several changes in the relative rankings of the classifiers. In some cases, the same classifier might figure among the best and the worst methods as we vary the k parameter in the narrow range of $[4, 14]$. These changes of relative order may cause some serious concerns for the research community. Without an extensive analysis of the influence of parameter values, some studies may experimentally show that a proposed algorithm outperforms the state-of-the-art algorithms, being that this conclusion only holds for certain parameter values. We are not claiming here that such an incident has ever happened, and we have not observed any such evidence, however; it is certainly undesirable for the research community to be affected of such a situation.

We suggest that every research paper that proposes a new SSL algorithm or graph construction method to fully analyze the influence of its parameters. The experimental setup used in this paper is a proposal of how newly

proposed methods should be evaluated. It is important to evaluate the algorithms' performance for a wide range of *external* parameters, such as k , and graph construction methods. Some algorithms also have *internal* parameters, such as regularization parameters, that also need to be evaluated, as we show next.

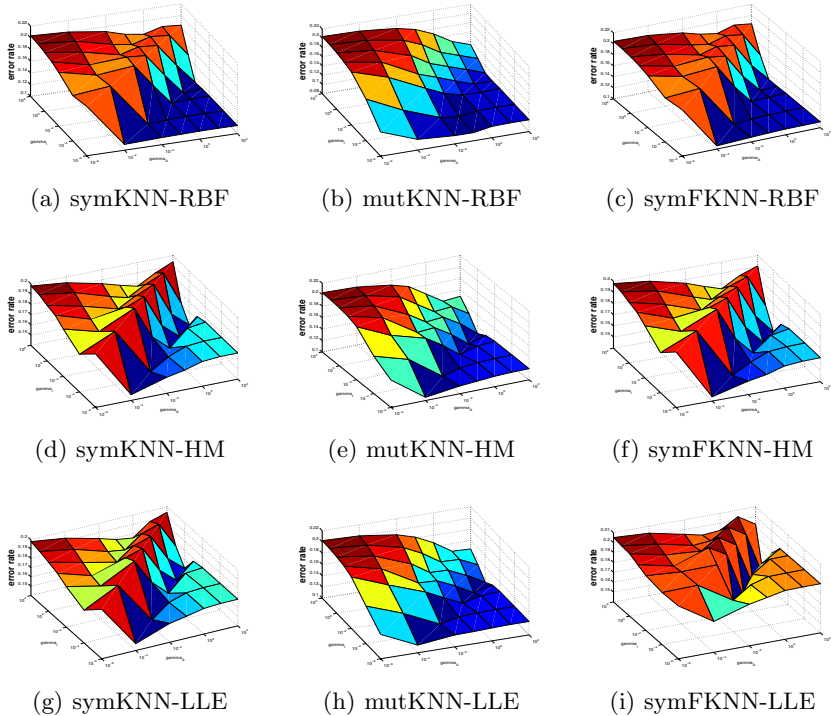


Fig. 3. Error surfaces for the LapRLS algorithm on the USPS data set

Influence of regularization parameters. We evaluate the influence of regularization parameters on the classification performance of the graph-based SSL algorithms. We evaluate the error surfaces generated by the SSL algorithms for each graph construction method and data set. Due to lack of space, we only show the most relevant results for LapRLS and LapSVM. We fixed the value of k that achieved the best error rate for each combination of SSL algorithm and graph construction method. In the sequence, we varied the values of γ_A and γ_I , as described in Section 4.3.

Fig. 3 shows the results for LapRLS on the USPS data set for different graph construction methods. We see that mutKNN generated smoother error surfaces than symKNN and symFKNN graphs, independent of the weighted matrix generation method used.

Many of the obtained results for this analysis are qualitatively equivalent fixing an SSL algorithm and a data set. However, we found some specific results that have an explicit pattern for parameters choice and others which may not have any evident pattern. We discuss them in the following.

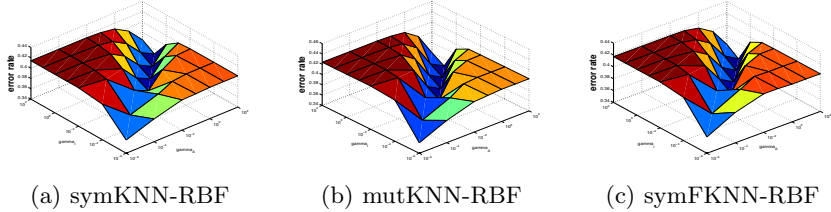


Fig. 4. Error surfaces for LapRLS on the TEXT data set using the RBF kernel

Fig. 4 shows the obtained results for LapRLS on the TEXT data set using the RBF kernel combined with the adjacency graph construction methods. We see that the “optimal region” occurs only when $\gamma_A = \gamma_I$.

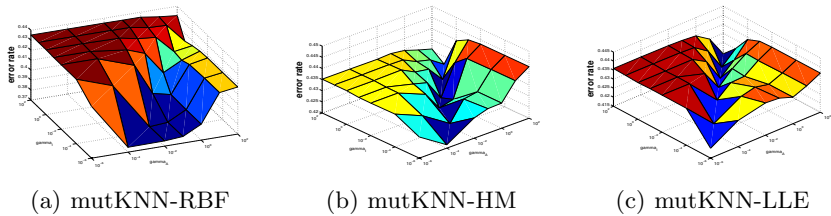


Fig. 5. Error surfaces for LapSVM on the TEXT data set using the mutKNN graph

Fig. 5 shows the obtained results for LapSVM on the TEXT data set using the mutKNN graph combined with the weighted matrix generation methods. We can not see any evident pattern that could help parameter choice. This may be an obstacle to apply LapSVM on real applications on text classification. For instance, the “optimal region” for the mutKNN-LLE graph occurs when $\gamma_A = \gamma_I$, which is not a good setting for the other graphs.

5 Conclusions and Further Research

In this paper, we provided a detailed empirical comparison of five state-of-the-art, graph-based SSL algorithms combined with three adjacency graph construction methods and three weighted matrix generation methods. Our experimental evaluation indicated that the SSL algorithms are strongly affected by the graph sparsification parameter value and the choice of the adjacency graph

construction and weighted matrix generation methods. The algorithms that have regularization parameters were also very dependent on a good setting of these parameters.

Consequently, we proposed an experimental setup that should be used in empirical comparisons in future work in SSL. Our idea is that a newly proposed algorithm should not be compared to other state-of-the-art algorithms using only the best case analysis. We believe that a detailed evaluation of all parameters is necessary. Due to the nature of SSL, in which there exists only a limited number of labeled examples, tuning all parameters might be unfeasible. Therefore, there is a need for algorithms that are slightly dependable on parameter tuning, i.e., that have a stable performance over the parameter space.

Our experimental results showed a superiority of mutKNN over the symKNN and symFKNN graphs. However, our results also showed that mutKNN is unstable for a relatively small value of k . In addition, we showed that mutKNN tends to generate smoother error surfaces than symKNN and symFKNN graphs. Our experiments also indicated a superiority of the RBF kernel in comparison to the HM and LLE methods.

As we analyzed our experimental results, we noticed other interesting patterns that we could not verify given the lack of experimental evidence. We propose an investigation concerning the validity of these observations as future research. Our empirical observations are as follows:

- Although RMGT achieved satisfactory results on most data sets, it did not perform well on the USPS data set. As USPS is an imbalanced dataset, a possible explanation is that RMGT is not effective on data sets with imbalanced labels;
- Maier et al. [10] have pointed out that the mutKNN graph should be chosen if one is only interested in identifying the “most significant” cluster. Based on this statement, we suppose that mutKNN is the best graph when we are dealing with data sets with imbalanced labels because it may identify the “most significant” class (the minority class in this case). This hypothesis is supported by the fact that, in Table 1, mutKNN achieved better classification performance than symKNN and symFKNN for all combinations of SSL algorithm and weighted matrix generation method on the USPS data set;
- Table 1 shows that RMGT achieved the best overall classification performance in 4 out of 6 data sets. This surprising classification performance may be due to the addition of the normalization constraints $\mathbf{F}\mathbf{1}_c = \mathbf{1}_n$ and $\mathbf{F}^T\mathbf{1}_n = n\omega$ in the optimization framework. It would be interesting to investigate if other SSL algorithms’ classification performances could be improved if these constraints were included in their optimization framework;
- In Fig. 4, we observed that the “optimal region” occurs only when $\gamma_A = \gamma_I$. Since this behavior occurred for all graphs (the other results are not shown here due to lack of space), we ask if this setting should be chosen for text classification tasks when using LapRLS.

Acknowledgments. This research was supported by the Brazilian agencies CAPES and FAPESP. Thanks to Diego F. Silva, Vinícius M. A. de Souza, Rafael Giusti, Ricardo M. Marcacini, and Rafael G. Rossi for their help in the experiments.

References

1. Belkin, M., Niyogi, P., Sindhvani, V.: Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *JMLR* 7, 2399–2434 (2006)
2. Beyer, K., Goldstein, J., Ramakrishnan, R., Shaft, U.: When is “nearest neighbor” meaningful? In: Beeri, C., Bruneman, P. (eds.) *ICDT 1999*. LNCS, vol. 1540, pp. 217–235. Springer, Heidelberg (1998)
3. Chapelle, O., Schölkopf, B., Zien, A.: *Semi-supervised learning*. The MIT Press (2006)
4. Demšar, J.: Statistical comparisons of classifiers over multiple data sets. *JMLR* 7, 1–30 (2006)
5. Hein, M., Maier, M.: Manifold denoising. In: *NIPS 19*, pp. 561–568 (2007)
6. Jebara, T., Wang, J., Chang, S.F.: Graph construction and b-matching for semi-supervised learning. In: *ICML*, pp. 441–448 (2009)
7. Johnson, R., Zhang, T.: On the effectiveness of laplacian normalization for graph semi-supervised learning. *JMLR* 8, 1489–1517 (2007)
8. Liu, W., Chang, S.F.: Robust multi-class transductive learning with graphs. In: *CVPR*, pp. 381–388 (2009)
9. Liu, W., He, J., Chang, S.F.: Large graph construction for scalable semi-supervised learning. In: *ICML*, pp. 679–686 (2010)
10. Maier, M., Hein, M., von Luxburg, U.: Optimal construction of k -nearest-neighbor graphs for identifying noisy clusters. *Theoretical Computer Science* 410(19), 1749–1764 (2009)
11. Melacci, S., Belkin, M.: Laplacian support vector machines trained in the primal. *JMLR* 12, 1149–1184 (2011)
12. Ozaki, K., Shimbo, M., Komachi, M., Matsumoto, Y.: Using the mutual k -nearest neighbor graphs for semi-supervised classification of natural language data. In: *CoNLL*, pp. 154–162 (2011)
13. Roweis, S., Saul, L.: Nonlinear dimensionality reduction by locally linear embedding. *Science* 290(5500), 2323–2326 (2000)
14. Zhou, D., Bousquet, O., Lal, T., Weston, J., Schölkopf, B.: Learning with local and global consistency. In: *NIPS 16*, pp. 321–328 (2004)
15. Zhu, X.: *Semi-supervised learning literature survey*. Tech. Rep. 1530, Computer Sciences, University of Wisconsin-Madison (2005)
16. Zhu, X., Ghahramani, Z., Lafferty, J.: Semi-supervised learning using gaussian fields and harmonic functions. In: *ICML*, pp. 912–919 (2003)