

On tuning a mean-field model for semi-supervised classification

Emílio Bergamim & Fabricio Breve

Avenida 24A, 1515, Rio Claro, State of São Paulo, Brazil

E-mail: emiliobergjr@gmail.com

February 2022

Abstract. Semi-supervised learning (SSL) has become an interesting research area due to its capacity for learning in scenarios where both labeled and unlabeled data are available. In this work, we focus on the task of transduction - when the objective is to label all data presented to the learner - with a mean-field approximation to the Potts model. Aiming at this particular task we study how classification results depend on β and find that the optimal phase depends highly on the amount of labeled data available. In the same study, we also observe that more stable classifications regarding small fluctuations in β are related to configurations of high probability and propose a tuning approach based on such observation. This method relies on a novel parameter γ and we then evaluate two different values of the said quantity in comparison with classical methods in the field. This evaluation is conducted by changing the amount of labeled data available and the number of nearest neighbors in the similarity graph. Empirical results show that the tuning method is effective and allows NMF to outperform other approaches in datasets with fewer classes. In addition, one of the chosen values for γ also leads to results that are more resilient to changes in the number of neighbors, which might be of interest to practitioners in the field of SSL.

1. Introduction

Semi-supervised learning (SSL) has become a major area of interest in machine learning due to its capacity for learning in the presence of labeled and unlabeled data. Conceptually, SSL is mainly seen as a midpoint between the major areas of *unsupervised* and *supervised* learning [1, 2].

In unsupervised tasks, the main goal is to determine the intrinsic structure of a dataset D , such as estimating a density function over it or partitioning it into subsets whose inner elements carry some form of similarity (also known as *clustering*). Supervised learning, on the other hand, aims to learn a function $y(x)$ from a sample of pairs $D = \{(x_i, y_i)\}_{i=1}^N$, like in the case of regression and classification [3].

The applicability of these two paradigms is usually defined according to the amount of labeled data available to the learner: unsupervised learning is suited for tasks where no labeled data is available *a priori*, while supervised learning is used when D is completely

labeled [4]. Such an extreme separation would lead to asking what to do when only a fraction of the elements of D are labeled. This is the scenario where *semi-supervised classification* takes place [2].

Under partially labeled data, one may be interested in learning a rule to label unseen data or in labeling all available data. The first situation is known as *inductive*, while the second is called *transductive* semi-supervised classification [2].

Algorithms for transduction generally rely on representing D through a graph for which each node represents an instance and weighted edges among these represent the similarity of a pair of elements of D [2, 5]. Transductive algorithms have found recent applications in many areas such as image segmentation [6], online tracking [7], text recognition in images [8] and video object segmentation [9].

1.1. Contributions

Our work revolves around semi-supervised transduction with the Potts model [10]. Previous works with similar models have pointed out the problem of tuning the β parameter as the main obstacle to the employment of these models both in clustering [11] and transduction [12, 13] tasks.

We then study how classification results depend on β by analyzing results in a range of the said parameter. In this investigation, we find that the optimal phase is highly dependent on the amount of labeled data available and that a more stable region for choosing β is associated with the probability of the most probable configuration of the system, Γ . We then construct an approximation for Γ and propose tuning β by choosing a target value for Γ .

Next, we compare the proposed approach with classical SSL approaches using two different targets for Γ . To do so, we investigate the behavior of said algorithms in different constructions of the similarity graphs and different rates of labeled data. We find our tuning approach is effective and can outperform existing approaches on some datasets, especially when the number of classes is small.

Also, one of the tuning methods introduced showed to be particularly resilient to changes in the number of neighbors in the graph. This property is particularly useful, since finding an optimal topology for the similarity graph is a difficult task.

1.2. Organization

The work follows with a brief exposure of transductive semi-supervised classification's particularities and the classical algorithms for this task. We then discuss the Potts model and the Naive Mean-Field Approximation that allows for the approximate calculation of statistical properties and connects this approach with other transductive ones through a propagation algorithm. Next, the experimental setup is introduced and followed by two different sections on experimental results, the first regarding the dependency and tuning of β and the second consisting of comparisons with classical algorithms.

2. Transductive semi-supervised classification

Let $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be a dataset for which $\mathbf{x}_i \in \mathbb{R}^n$ denotes the attributes of the i -th sample and $y_i \in \{1, \dots, q\}$ its label. In a SSL setup there is a subset $D_l \subset D$ for which the labels are known in advance and another subset $D_u = D \setminus D_l$ of unlabeled samples.

For transductive methods, the goal is to label all elements of D . To do so, one first constructs a *weighted similarity graph* G_D for which each node is mapped to a sample of D and each edge describes the similarity among neighboring samples. Next, one applies an inference algorithm to the graph to label each node. These are the two main steps of semi-supervised transduction [2, 5], which we will now discuss in more depth.

2.1. Transductive inference in graphs

The goal of a transductive algorithm is to use the information available in D_l to label points in D_u [5]. To do so, one needs a similarity matrix \mathbf{W} for which entries $W_{i,j} \in [0, 1]$ describe the similarities between the i -th and j -th samples in D . The nonzero entries in \mathbf{W} correspond to the (weighted) edges of G_D .

Once such objects are available to the learner, the inference phase consists of a propagation algorithm where the *a priori* knowledge in D_l is propagated over G_D (and weighted by \mathbf{W}) to all samples in D . In this context, two algorithms are currently established as the standard approaches: *Gaussian Random Fields* (GRF) and *Local and Global Consistency* (LGC) [2, 5], the first being regarded as the state of the art by some authors [16, 17].

If there are q possible labels in a set, let $\phi_i(s_i)$ be the probability that the i -th sample in D has label $s_i \in \{1, \dots, q\}$. GRF [18, 19] is an approach that minimizes the objective function

$$H_{GRF}(\phi) = -\frac{1}{2} \sum_{i,j} W_{i,j} \sum_s [\phi_i(s) - \phi_j(s)]^2 \quad (1)$$

under the constraint that instances in D_l have their probabilities sampled to

$$\phi_i(s_i) = \begin{cases} 1, & \text{if } s_i = y_i \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

LGC [20], on the other hand, was proposed as a way to overcome limitations of GRF regarding noisy labels or irregular graph structures and is also widely used [2, 5]. In this case, one does not work with probabilities in the strict sense, but with a set of vectors $\{\mathbf{f}_i\}_{i=1}^N$ in \mathbf{R}^q that minimizes

$$H_{LGC}(\mathbf{f}) = \sum_{i,j} W_{i,j} \|\mathbf{f}_i - \mathbf{f}_j\|^2 + \frac{1-\alpha}{\alpha} \sum_i \|\mathbf{f}_i - \boldsymbol{\theta}_i\|^2, \quad (3)$$

where $\boldsymbol{\theta}_i \in \mathbf{R}^q$ is defined by

$$\theta_{i,s} = \begin{cases} 1, & \text{if } (x_i, y_i) \in D_l \text{ and } s = y_i, \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

and $\alpha \in (0, 1)$ is the parameter of such model [20].

An exact approach to minimize the objectives in (1) and (4) requires the inversion of \mathbf{W} , for which complexity is $O(N^3)$ [2, 20]. However, both methods allow for an iterative solution with complexity $O(N^2)$. These are presented in algorithms 2.1 and 2.2. After convergence of such procedures the points are labeled as

$$y_i = \operatorname{argmax}_{s_i} \phi_i(s_i) \quad \text{and} \quad y_i = \operatorname{argmax}_{s_i} f_{i,s_i} \quad (5)$$

for GRF and LGC, respectively.

Algorithm 2.1 Iterative GRF

Input: $\phi^{(0)}$, \mathbf{W} , D_l , ϵ , t_{max}
 $t \leftarrow 0$, $\delta \leftarrow \epsilon$
while $t < t_{max}$ and $\delta \geq \epsilon$ **do**
 $\delta \leftarrow 0$
 for $i = 1, \dots, N$ **do**
 if $(x_i, y_i) \notin D_l$ **then**
 for $s = 1, \dots, q$ **do**
 $\phi_i^{(t+1)}(s) \leftarrow \frac{\sum_{j \neq i} W_{i,j} \phi_j^{(t)}(s)}{\sum_l \sum_{j \neq i} W_{i,j} \phi_j^{(t)}(l)}$
 end for
 end if
 $\delta \leftarrow \max\{\delta, \max_s |\phi_i^{(t+1)}(s) - \phi_i^{(t)}(s)|\}$
 end for
 $t \leftarrow t + 1$
end while
Output: $\phi^{(t)}$

Algorithm 2.2 Iterative LGC

Input: $\mathbf{f}^{(0)}$, \mathbf{W} , $\boldsymbol{\theta}$, α , ϵ , t_{max}
 $t \leftarrow 0$, $\delta \leftarrow \epsilon$
while $t < t_{max}$ and $\delta \geq \epsilon$ **do**
 $\delta \leftarrow 0$
 for $i = 1, \dots, N$ **do**
 for $s = 1, \dots, q$ **do**
 $f_{i,s}^{(t+1)} \leftarrow \alpha \sum_{j \neq i} W_{i,j} f_{j,s}^{(t)} + (1 - \alpha) \theta_{i,s}$
 end for
 $\delta \leftarrow \max\{\delta, \max_s |f_{i,s}^{(t+1)} - f_{i,s}^{(t)}|\}$
 end for
 $t \leftarrow t + 1$
end while
Output: $\mathbf{f}^{(t)}$

2.2. Similarity construction

Both algorithms presented before are examples of propagation dynamics on graphs. In fact, by examining these methods, one can see that the weighting process becomes irrelevant if G_D describes the relations among elements of D in such a way that edges only connect points that belong in the same class. Weighting becomes a necessity since one usually does not know in advance a sufficient topology of G_D for class detection.

In this work a k nearest neighbors (kNN) approach is used for construction of \mathbf{W} : the nonzero entries of the i -th row of \mathbf{W} are the k nearest neighbors of x_i according to a distance function $d(\cdot, \cdot)$.

A highly popular method for weighting edges uses a gaussian radial basis function [2]

$$W_{i,j} = \exp \left\{ \frac{d(\mathbf{x}_i, \mathbf{x}_j)^2}{2\sigma^2} \right\} \quad (6)$$

and the parameter σ can be tuned as [21]

$$\sigma = \frac{1}{3N} \sum_i d_{i,k(i)}, \quad (7)$$

where $k(i)$ denotes the k -th nearest neighbor of the i -th element of D . A previous study using several weighting schemes reported this approach as the most effective when combined with different inference algorithms [22].

It is also a common place in the literature to use a sparse construction of \mathbf{W} [22, 2, 5] where most of its entries are set to zero, leaving only the most significant entries in order to increase contrast among different classes. Choosing how sparse one wants \mathbf{W} to be is the problem of choosing a particular k to the problem.

In [22] it was verified that a robust way of obtaining a sparse similarity is by setting the non-zero entries of \mathbf{W} to initially be the k -th nearest neighbor of each sample in D and tune σ as in Equation 7. This is then followed by a symmetrization and sparsification procedure from which a novel similarity is obtained via

$$W_{i,j} \leftarrow \min\{W_{i,j}, W_{j,i}\}. \quad (8)$$

Then, \mathbf{W} is normalized by setting the sum of its rows to 1:

$$W_{i,j} \leftarrow \frac{W_{i,j}}{\sum_{l \neq i} W_{i,l}}. \quad (9)$$

In the present work we will focus on the above similarity construction with the distance function $d(\cdot, \cdot)$ being the euclidean (L2) norm and evaluate algorithms for different values of k . For further reading on the problem of constructing \mathbf{W} we refer the reader to [22, 23, 2, 5].

3. Potts model and naive mean fields

The main hypothesis of semi-supervised learning is that elements of D with the same label belong to the same cluster - i.e., the classes of the dataset are the labeling of different clusters [1, 2].

Under this consideration, transduction can be viewed as the optimization of a cost function $H_C(\mathbf{s})$ over a set of labels $\mathbf{s} = (s_1, \dots, s_N)$ under the restriction that known labels are constants. In fact, the case

$$H_C = - \sum_{i < j} W_{i,j} \delta(s_i, s_j), \quad (10)$$

where $\delta(\cdot, \cdot)$ is the Kronecker Delta, is known to be equivalent to the GRF algorithm when variables corresponding to labeled data are held fixed [24].

Equation 10 is known as the Potts model [10] without an external field. In other works, this model has been used in the presence of an external field as defined for the LGC model in Equation 4 [12, 13]. In this case, the cost function is of the form

$$H(\mathbf{s}) = - \sum_{i,j} W_{i,j} \delta(s_i, s_j) - \sum_{i,s} \theta_{i,s} \delta(s_i, s). \quad (11)$$

Also, the latter approach is probabilistic since labels are attributed according to the distribution

$$\Psi(\mathbf{s}) = \frac{1}{Z_H} \exp\{-\beta H(\mathbf{s})\}, \quad (12)$$

where Z_H is a normalization constant known as the partition function, for which calculation imposes a problem, as its exact evaluation demands one to sum over all q^N possible configurations of \mathbf{s} .

A possible overcome of this problem is the usage of Monte Carlo methods, as was done both in clustering (corresponding to the zero-field form of Equation 11) and transduction [14, 15] applications. This, however, has the problem of relying on random number generation, which renders some unpredictability to the behavior of the algorithms.

A deterministic approach to the problem is the use of mean-field methods, more recently studied both in clustering [25] and transduction [12, 13]. In this case, one seeks to approximate Equation 12 by a more tractable distribution, with the easier way of doing so being known as the Naive Mean Field (NMF) method [26, 12].

This approach consists of finding a distribution $\Phi(\mathbf{s})$ that minimizes the Kullback-Leibler divergence

$$D_{KL}(\Phi||\Psi) = \sum_{\mathbf{s}} \Phi(\mathbf{s}) \ln \left\{ \frac{\Phi(\mathbf{s})}{\Psi(\mathbf{s})} \right\} \quad (13)$$

under the constraint that Φ is a product distribution:

$$\Phi(\mathbf{s}) = \prod_i \phi_i(s_i). \quad (14)$$

Under this construction Φ is a distribution that considers \mathbf{s} as a vector of independent variables. Therefore, $\phi_i(s_i)$ is the marginal distribution of s_i . Minimization of Equation 13 with respect to such marginals results in a set of non-linear coupled equations

$$\phi_i(s_i) = \frac{1}{Z_i} \exp\{h_i(s_i)\}, \quad \text{with } Z_i = \sum_{s_i} \exp\{h_i(s_i)\} \quad (15)$$

and

$$h_i(s_i) = \beta \left(\theta_{i,s_i} + \sum_{j \neq i} W_{i,j} \phi_j(s_i) \right) \quad (16)$$

that can be solved iteratively with complexity $O(N^2)$ as shown in Algorithm 3.1. After this procedure, one labels instances of D in a similar way as is done for GRF and LGC:

$$y_i = \operatorname{argmax}_{s_i} \phi_i(s_i). \quad (17)$$

Algorithm 3.1 Iterative NMF

Input: $\Phi^{(0)}$, β , \mathbf{W} , $\boldsymbol{\theta}$, ϵ , t_{max}
 Initialize $\mathbf{h}^{(0)} = \mathbf{0}$
 $t \leftarrow 0$, $\delta \leftarrow \epsilon$
while $t < t_{max}$ and $\delta \geq \epsilon$ **do**
 $\delta \leftarrow 0$
 for $i = 1, \dots, N$ **do**
 for $s = 1, \dots, q$ **do**
 $h_i(s)^{(t+1)} \leftarrow \beta \left(\theta_{i,s_i} + \sum_{j \neq i} W_{i,j} \phi_{j,s}^{(t)} \right)$
 end for
 Calculate $\phi_i^{(t+1)}$ using Equation 15.
 $\delta \leftarrow \max\{\delta, \max_s |h_i(s)^{(t+1)} - h_i(s)^{(t)}|\}$
 end for
 $t \leftarrow t + 1$
end while
Output: $\Phi^{(t)}$

3.1. Previous works on tuning β

We will work with the NMF approximation to calculate the statistical properties of the Potts model, the remaining problem is to understand how the classification results will be affected by β .

Earlier works on clustering (corresponding to the zero-field form Equation 11) advocate for the existence of a range for β where results are optimal in some sense [14]. It is known, however, that in such a range different clustering structures exist and the idea of optimality can only emerge by analyzing and combining these different possible clusterings [11].

The application of these ideas to classification is followed immediately by fixing variables associated with elements of D_l to known labels and applying the same tuning procedure to estimate β . It was also observed that increasing the size of D_l also increases the optimal range for β [15].

A major drawback of the above works is the necessity for evaluating the statistical properties of the model for different values of β in order to determine the optimal range.

Despite using Monte Carlo methods, the usage of mean-field methods would not offer a significant improvement to this, especially in the semi-supervised classification problem, since GRF is a well-established non-parametric approach.

Regarding mean-field methods, more recent work was done with the Ising model [12, 13] using a similar similarity construction as the one outlined in the previous section. However, the authors in these works highlighted that the need for efficient tuning of β is the main obstacle to practical usage, since results on accuracy point that this approach reaches state-of-the-art performance [12, 13].

It is also noteworthy that earlier work on SSL with the Potts model uses β -independent fields set to be infinitely strong [15], while most recent approaches [12, 13] have focused on β -dependent fields, as is the case of the model in Equation 12. Both approaches do so in order to keep labeled data fixed to their known labels, but in different ways. The first approach freezes variables to their known labels. The second allows marginals of variables with non-zero fields to change with β , but not their labels. To better understand this case, we draw attention to the normalization of similarities (Equation 9) and the fact that marginals are upper-bounded by unity to obtain the following inequalities for the NMF equations (Equation 16)

$$\beta\theta_{i,s_i} \leq h_i(s_i) \leq \beta(\theta_{i,s_i} + 1). \quad (18)$$

Then, the definition of fields (Equation 4) implies that if the m -th instance of D is labeled as y_m we have

$$h_m(s_m) \leq h_m(y_m). \quad (19)$$

We then note that equality in the above expression can only be achieved for $s_m \neq y_m$ and $\beta \neq 0$ if

$$\sum_{j \neq m} W_{m,j}(\phi_j(s_m) - \phi_j(y_m)) = 1 \quad (20)$$

and normalization of similarities implies this can only be achieved if $\phi_j(s_m) = 1$ for all neighbors of m . This would imply marginals of neighbors of m to be unaffected by ϕ_m , which is true only if we have $W_{m,j} = 0$ together with $W_{j,m} = 0$ which is not impossible under the similarity construction discussed in subsection 2.2. We then conclude that algorithm 3.1 is not able to change known labels in D .

Finally, we also highlight recent work on clustering using a Potts spin-glass together with mean-field methods [25], which related optimal results to a phase transition. Our work, however, will focus on ferromagnetic interactions to follow along previous works on semi-supervised classification with mean-field methods that were proven effective in this task [12, 13].

Next, we discuss the experimental setup that will be used for the remainder of this paper in order to understand the β dependency problem and compare NMF, LP and GRF.

4. Experimental setup

As we aim to evaluate semi-supervised methods for different configurations of the problem, here we discuss the data that will be used for such evaluation. Particularly, three bidimensional artificially generated datasets will be used together with other six high-dimensional datasets available in the literature.

The bidimensional datasets (Figure 1) are **Two moons**, a set consisting of two balanced non-convex classes; **Three clusters**, a set of three unbalanced classes, with two of them being non-convex and **Five Gaussians** with different locations and standard deviations, as well as different proportions. The first and last datasets contain 1000 elements, while the second contains 900.

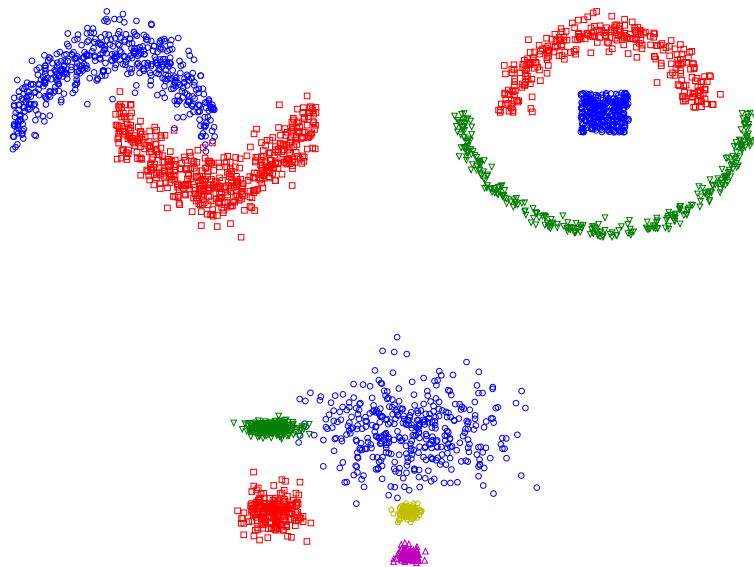


Figure 1. Images of the three bidimensional datasets used for benchmarking: two moons (top left), three clusters (top right) and five gaussians (bottom). Different symbols relate to different labels.

The high-dimensional datasets are

- **Digit1** was introduced in [1] to evaluate semi-supervised algorithms. It has $N = 1500$, dimension $d = 241$ and was preprocessed by the original authors. It is a set of artificial images of the digit “1”, divided in two classes [1].
- **Twonorm** and **Ringnorm** are available in the DELVE repository [27]. These are sets of 7400 instances in a 20-dimensional space consisting of two gaussian distributions. In each dataset classes are differed by their means and covariance matrices.
- **Landsat** is a dataset of 6435 3×3 hyperspectral satellite images divided in 7 classes. It was obtained from the UCI repository [28].

- **USPS** is a set of handwritten digits from “0” to “9”, but with 9298 images [29] and dimensions 16×16 .
- **Texture** is a set from the ELENA project consisting of 5500 patterns describing 11 different textures (the classes). Each instance is described by 40 attributes estimated from fourth-order modified moments. As the original repository is missing, we refer the reader to the version in [30].

We evaluate algorithms over different constructions of D_l . Letting N_l be the number of elements in D_l and $r_l = N_l/N$ the rate of labeled data presented to the learner, for each value in the grid $r_l \in [0.02 : 0.2 : 0.02]$ we randomly generate twenty different realizations of D_l assuring that at least one instance of each class is presented.

Classification results are then evaluated according to the accuracy and adjusted mutual information (AMI, evaluated in *nats*) [31]. The second is a metric commonly used to evaluate different clusterings of a dataset. Due to the clustering hypothesis of semi-supervised learning, comparing the accuracy and AMI will allow us to have a better understanding of how clustering connects to classification in these models.

To run experiments on the described datasets we implement the propagation algorithms 2.1, 2.2 and 3.1 in the C programming language [32] and then build an interface in Cython [33] to make these functions callable via Python [34].

The Python part of our experiments is mainly the handling of input and output for each dataset, while algorithms 2.1, 2.2 and 3.1 are executed in C using multithreaded parallelism via the OpenMP standard [35]. All codes used for benchmarking are available at <https://github.com/boureau93/ssl-nmf>. Experiments were carried out on an Intel i5-1135G7 processor with 8GB RAM.

5. Dependency of classifications on β

Our first experimental evaluation aims to study how classification results depend on β . As we wish to develop a tuning procedure for the said parameter, we must relate said results to a statistical property. In this work we will focus on the mode probability

$$\Gamma = \prod_i \phi_i(y_i) \quad (21)$$

and how this quantity connects to the accuracy and AMI at different values of β .

Figures 2 through 6 show results for accuracy, AMI, execution time and Γ as a function of β for three different values of r_l and a fixed topology of G_D at $k = \log_2 N$. In these experiments we also set $t_{max} = 10^4$ and $\epsilon = 10^{-3}$ and change β in the grid $10^{[-3.3:0.2]}$.

Experiments on the artificial bidimensional datasets (Figures 2-4) show that optimal classification results are associated with higher values of Γ . When $\Gamma = (\frac{1}{q})^N$ the model is in an equiprobable configuration, leading to the worst classification results since one cannot effectively distinguish different labels via their probabilities. The increase of β then leads to the increase of Γ . We also observed that optimal results are usually

associated with Γ above the midpoint between equiprobable probabilities and maximum probability, i.e., $\Gamma \geq (\frac{1+q}{2q})^N$.

As our experiments were done using $k = \log_2 N$ (Figures 2-6), we expect most edges to connect only nodes that belong to the same class, but some inter-class connectivity still exists. This is overcome by the increase in β , which increases correlations among variables and, together with the information given by labeled data increases contrast among classes by making wrongly connected nodes less relevant.

However, the existence of a range for β where accuracy and AMI are stable at a maximum seems to be conditioned on the amount of labeled data available. On the two moons (Figure 3) dataset we see this behavior: for $r_l = 0.02$ there is a peak of optimal classification that is followed by a slow decrease in accuracy and AMI, while for the higher values of r_l the model reaches an optimal plateau of such quantities.

On the two high-dimensional real datasets evaluated (Figure 5 and Figure 6) we see the absence of the optimal classification plateau for all evaluated values of r_l . The Landsat dataset shows a more stable optimal region than USPS, which has a more pronounced peak in accuracy and AMI.

This behavior is in line with previous studies [15] pointing that the addition of labeled data increases the range of β containing optimal results. On top of that, our results show that the existence of an optimal plateau where classifications are not affected by changes in β is conditioned on the size of D_l , while the sufficient amount for the said phenomenon to occur is dataset dependent. We can then conclude that increasing r_l increases the stability of classification metrics once Γ is big enough. As labeled data becomes scarcer, the difference between optimal results and results in the Γ plateau becomes more evident.

It is also notable that at higher Γ the execution time of the algorithm scatters to a plateau of higher execution that is up to three orders of magnitude slower than the low Γ region. Together with the decrease in accuracy and AMI at the Γ plateau shown by some datasets we conclude this is a region where elements at the boundary of a class can be difficult to label due to their high correlations to elements of different classes, demanding more iterations of algorithm 3.1. This can be evidenced by looking at the three clusters dataset, in which classes are more separated (Figure 1) and the execution time does not present a plateau at higher Γ .

The above also indicates that better classifications are not strictly associated with a higher execution time of the algorithm. What is true, however, is that higher values of Γ are associated with a decrease in computational performance that may or may not be in the form of a plateau.

In fact, for sufficiently high r_l , close to optimal values of accuracy and AMI can coexist with lower execution times as illustrated in the results for bidimensional datasets (Figures 2-4). Therefore, for an appropriate choice of β , increasing the amount of labeled data can lead to better computational performance.

Still regarding execution time, at higher values of β the algorithm can slow down up to three orders of magnitude when compared to the paramagnetic phase where $\Gamma \approx (\frac{1}{q})^N$.

We also highlight that, for higher values of β the used hardware lacks numerical precision, and solutions of the NMF equation return NaN values, leaving us unable to calculate Γ . So, the decrease in accuracy, AMI, and execution at higher β may not be associated with the model or the approximation, but due to an experimental limitation.

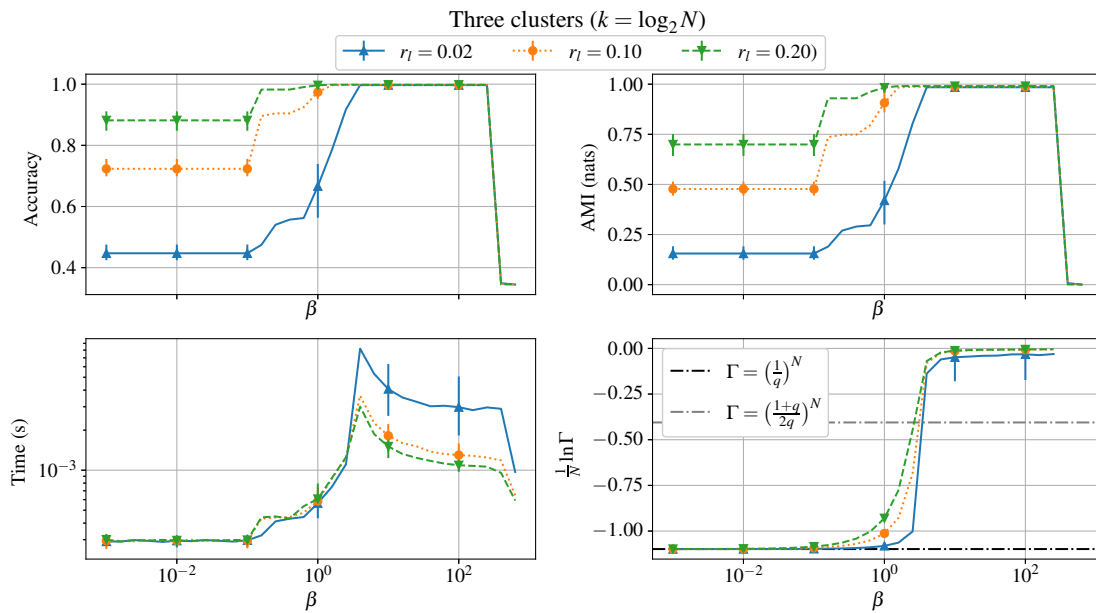


Figure 2. Results for the three clusters dataset as a function of β . Lines denote averages and bars denote maximum and minimum over different realizations of D_I .

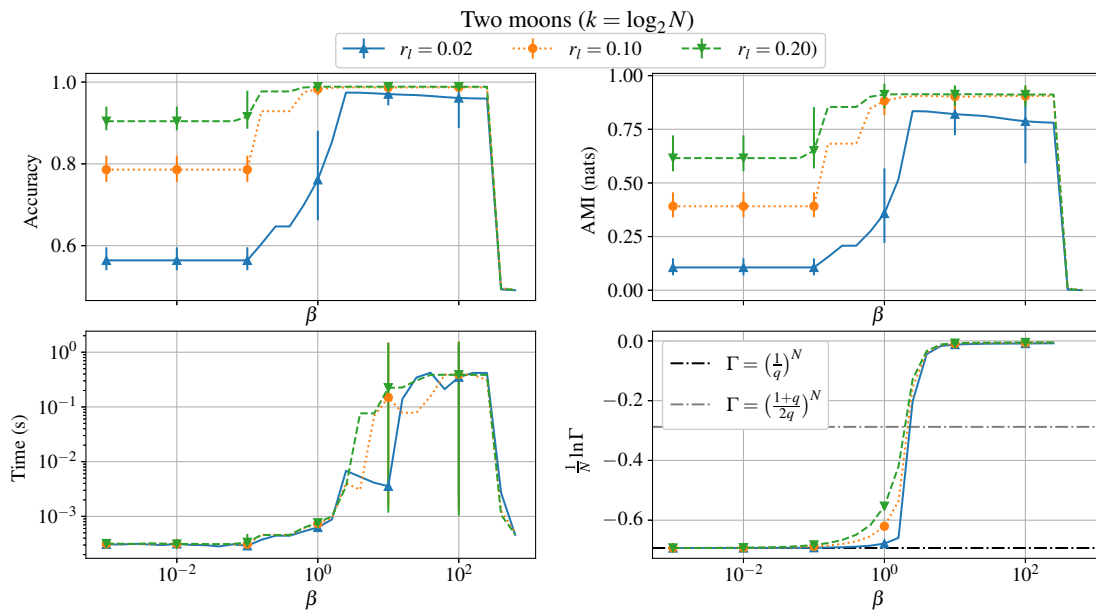


Figure 3. Results for the two moons dataset as a function of β . Lines denote averages and bars denote maximum and minimum over different realizations of D_I .

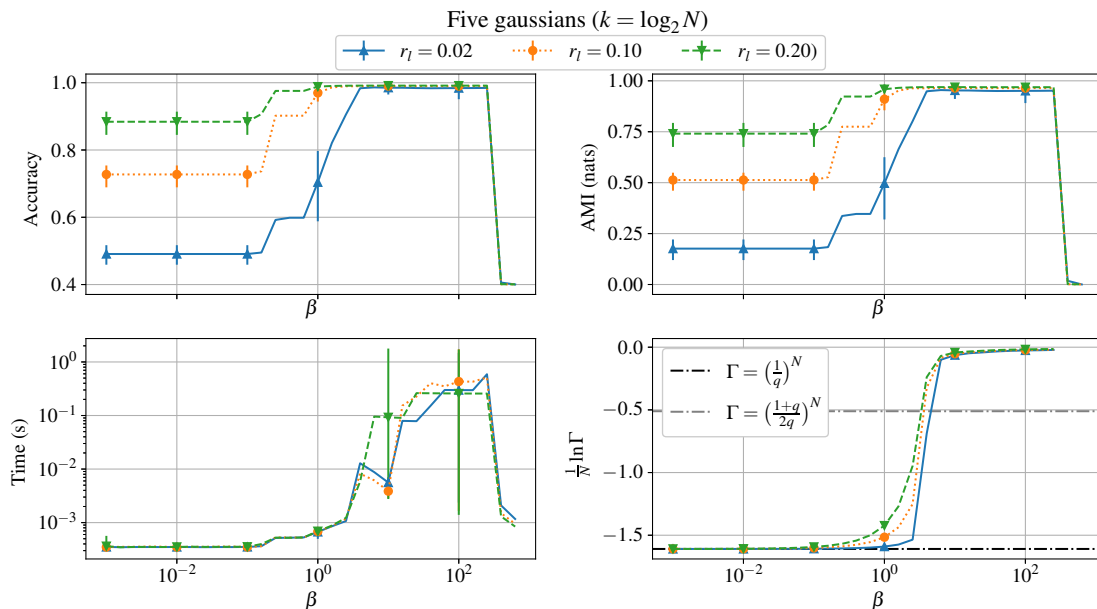


Figure 4. Results for the five gaussians dataset as a function of β . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

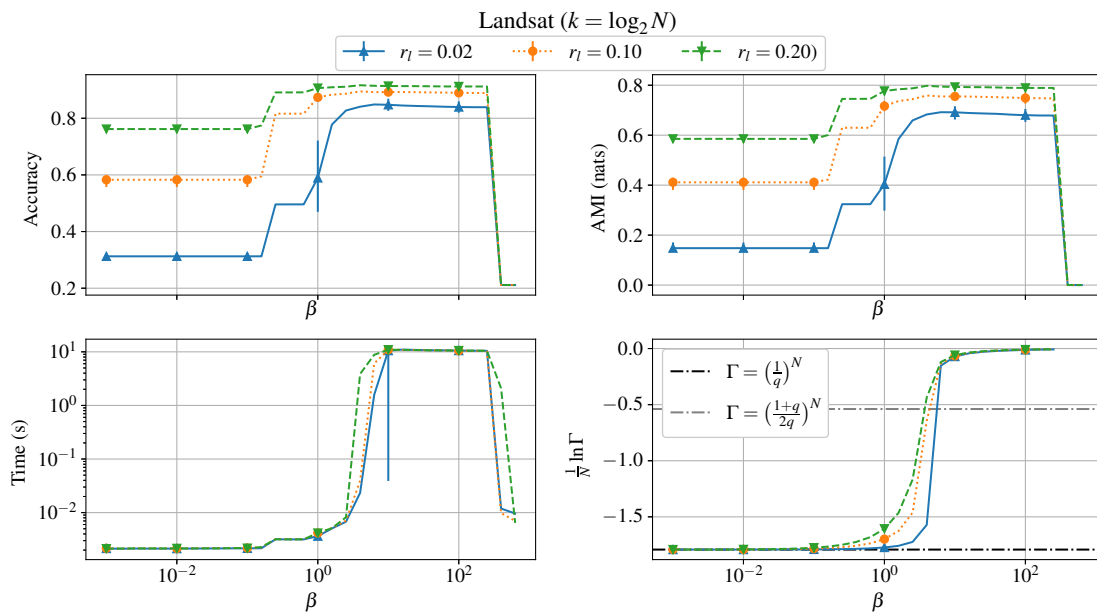


Figure 5. Results for the Phoneme dataset as a function of β . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

5.1. A tuning procedure for β

The discussion made so far points out that finding the pointwise optimal value of β is hard, specially since Γ does not show a profile that is similar to the the accuracy and AMI curves, therefore suboptimal configurations exist at higher Γ , but those seem to be probabilistically indistinguishable from optimal ones. However, the existence of a range

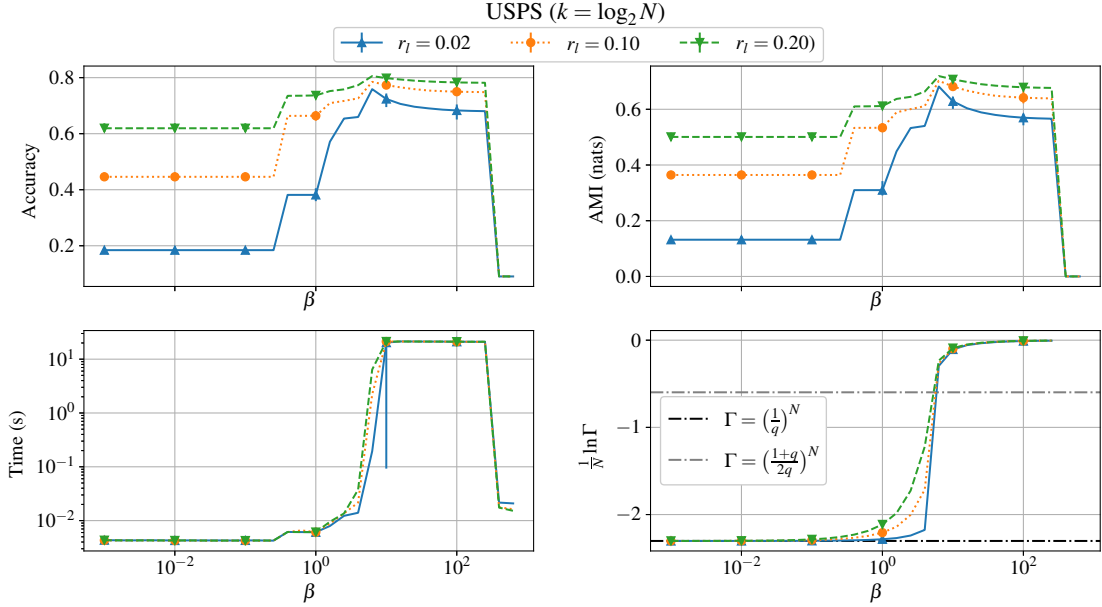


Figure 6. Results for the USPS dataset as a function of β . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

where classification metrics are more stable and the observation that this is associated with a plateau in Γ indicates that one could try to tune β by tuning Γ .

The above then motivates us to find an approximation for Γ to tune the model. We will work with the logarithm of Γ

$$\ln \Gamma = \sum_i (h_{i,y_i} - \ln Z_i) \quad (22)$$

in order to make calculations easier. As noted, for low β we have $\psi_i(s_i) \approx \frac{1}{q}$. At this regime we can also approximate Z_i by the first terms in the Taylor series of the exponential:

$$Z_i \approx \sum_{s_i} (1 + h_{i,s_i}) = q + \beta \left(\sum_{s_i} \theta_{i,s_i} + 1 \right), \quad (23)$$

where the last expression comes from the normalization of the rows of \mathbf{W} (Equation 9). Then,

$$\ln \Gamma \approx \sum_i \left[\beta \left(\theta_{i,y_i} + \frac{1}{q} \right) - \ln \left\{ q + \beta \left(\sum_{s_i} \theta_{i,s_i} + 1 \right) \right\} \right]. \quad (24)$$

Now, to evaluate the sum over all variables in the above expression we recall the definition of $\boldsymbol{\theta}$ in Equation 4. Since elements in D_l cannot change labels through NMF, summing all θ_{i,y_i} yields the number of elements in D_l . By the definition in Equation 4 we also note that the sum inside the logarithm is equal to 1 for labeled variables and 0 for unlabeled ones. Therefore, Equation 24 can be written as

$$\ln \Gamma \approx \beta \left(N_l + \frac{N}{q} \right) - N_l \ln \{ q + 2\beta \} - (N - N_l) \ln \{ q + \beta \} \quad (25)$$

that when divided by N yields

$$\frac{1}{N} \ln \Gamma \approx \beta \left(r_l + \frac{1}{q} \right) - r_l \ln\{q + 2\beta\} - (1 - r_l) \ln\{q + \beta\}. \quad (26)$$

It is noteworthy that this can be used only in the low β regime since it diverges to infinity as $\beta \rightarrow \infty$.

As we have shown before based on the empirical study presented in this section, more stable results of classification are located above a threshold in Γ . Therefore we aim to choose β by solving

$$\frac{1}{N} \ln \Gamma = \ln \gamma \quad (27)$$

for $\gamma \in (0, 1]$, where γ is an user-defined parameter. One could argue we are now simply changing the problem of choosing β to the problem of choosing γ , which is true. Choosing γ , however, points to choosing how sure one wants to be about a classification, with the true level of belief being given by the actual value of Γ . This does not indicate that a higher γ will leads to better classification results, but, together with the aforementioned stability at higher Γ , this procedure has a foundation for its applicability.

Figure 7 shows the solution β_γ^* of Equation 27 using the approximation in Equation 26 and setting $\gamma = \frac{1+q}{2q}$ and $\gamma = 1$ as target values. To solve the non-linear equation we used the implementation of Newton's method in the SciPy library [36] with initial guess set to unity and tolerance set to 10^{-3} .

As Equation 26 is designed to work at lower values of β , one may face a problem in its correctness in situations where the number of classes q is too high or when the amount of labeled data provided is too low, since these conditions increase the values of β_γ^* (Figure 7).

Now, the existence of a procedure for tuning β will allow us to make a more in-depth comparison of NMF with GRF and LGC. In the next section we will use both values of γ to calculate β_γ^* and compare these approaches with the classical approaches in the field of semi-supervised learning.

6. Comparative evaluation of SSL algorithms

In this section we evaluate GRF, LGC and NMF using the previously constructed tuning approach for β with $\gamma = \frac{1+q}{2q}$ and $\gamma = 1$. The first is the observed lower bound reported in the previous section, while the second should be the limit for γ where our approximation holds. For all algorithms we set $t_{max} = 10^4$ and $\epsilon = 10^{-3}$, while for LGC we set $\alpha = 0.99$ as in the original paper [20]. Experiments are then conducted for fixed values of k and r_l while varying the other quantity.

6.1. Bidimensional datasets

As shown in figures 8-10 the behavior of the three algorithms is highly determined by the datasets, which makes it difficult to point out the most effective approach.

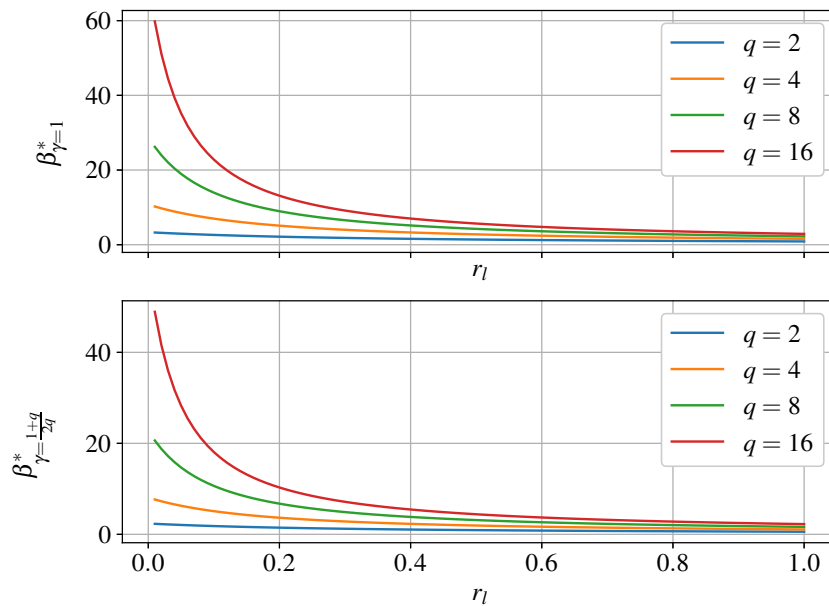


Figure 7. Solution of Equation 27 as a function of r_l for two different values of γ and four different values of q .

In the two moons dataset (Figure 8) setting $\gamma = \frac{1+q}{2q}$ provided a faster execution time at the expense of accuracy and AMI, especially in the regime of small r_l . The approach for $\gamma = 1$ also showed a better computational performance than GRF and LGC in most studied scenarios, but without the cost in classification metrics shown by the other tuning approach.

When looking at accuracy and AMI as functions of k in Figure 8 one sees that the difference in the two ways of choosing β is related to the first being less susceptible to increases in k . The second method has similar behavior to GRF and LGC in terms of classification, while being faster than both of these methods, especially for lower k .

The case of three clusters (Figure 9) showed to favor NMF approaches by a small portion in the low r_l regime. However, in denser topologies, GRF and LGC become faster, but the second shows a decrease in classification metrics for $k > 11$.

NMF shows an irregular execution time profile as a function of r_l in the five gaussians dataset, as illustrated in Figure 10. Since the procedure for choosing β depends strictly on r_l but not on the particular realizations of D_l , we see that the classification metrics of NMF can be highly affected by different configurations of the previously labeled instances of the dataset, particularly as labeled data becomes more scarce.

What we note as the similarity between the three studied datasets is that the increase in k tends to, on average, improve the execution time of the algorithms for $k > \log_2 N$, with the exception of LGC in the five gaussians dataset (Figure 10). This is quite interesting behavior as one could expect that the addition of edges to

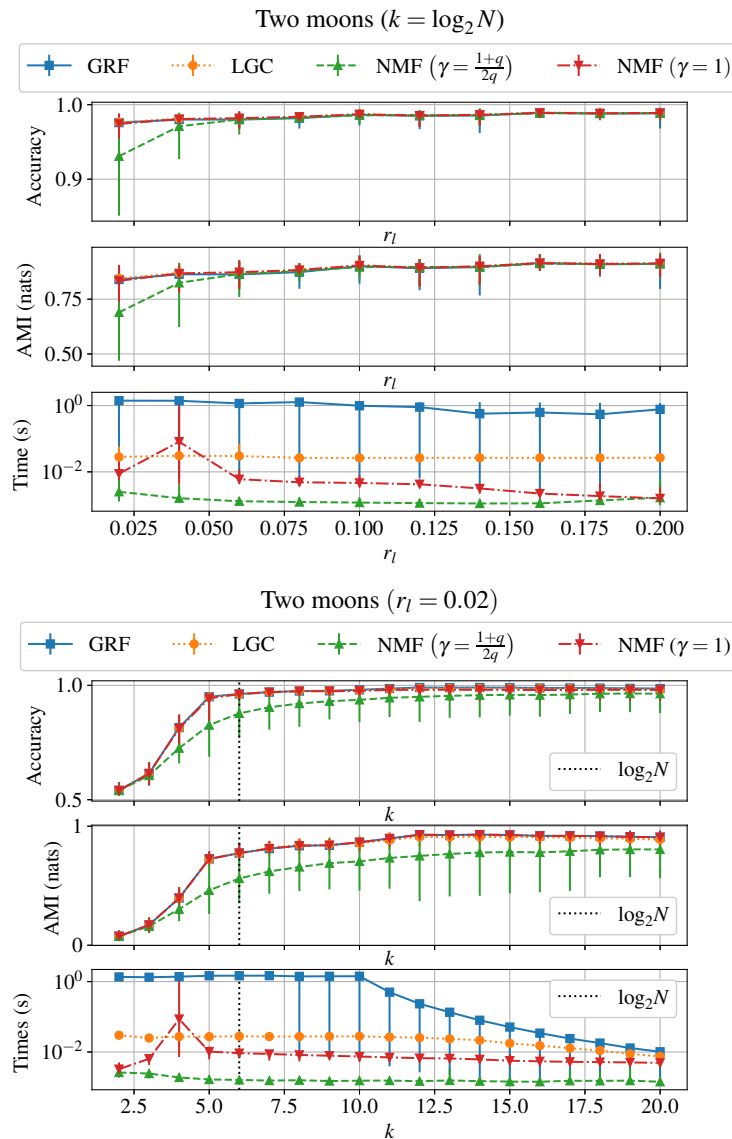


Figure 8. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

G_D would slow down these algorithms. We believe this is related to the procedure of symmetrization and sparsification discussed in subsection 2.2, which allows for the increase in k to connect more similar points and increase intra-class connectivity, making the problem easier to handle for the algorithms and enhancing their convergence.

The above is supported by the average improvement in classification metrics as k increases. The exception is again LGC, but in the three clusters dataset (Figure 9). Therefore our study so far leads us to believe that LGC (with $\alpha = 0.99$) is affected by changes in the topology of G_D in a different manner than GRF and NMF, which is an expected behavior based on a previous study that connects the Potts model with GRF [24].

We also highlight that in the bidimensional datasets (Figures 8-10) the accuracy

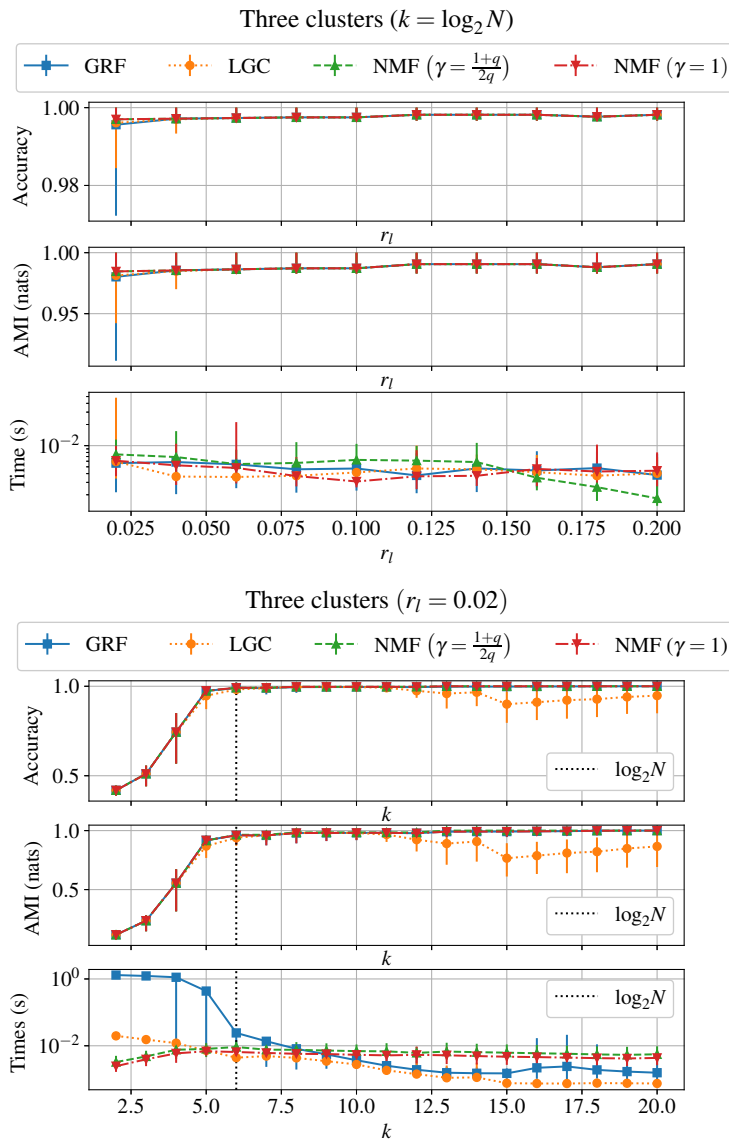


Figure 9. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

and AMI curves show a high agreement: they respond similarly to variations in r_l and k . Therefore, classification and clustering using GRF, LGC and NMF are closely related in the said datasets.

6.2. High-dimensional datasets

As we now move to analyze high-dimensional datasets we will also investigate data that comes from real-world phenomena, such as Landsat, Texture and USPS. Our discussion will navigate by each of the six sets being analyzed.

For Digit1 (Figure 11) we observe the utilization of NMF with $\gamma = 1$ yields the best results in accuracy and AMI when the problem is looked at as a function of r_l .

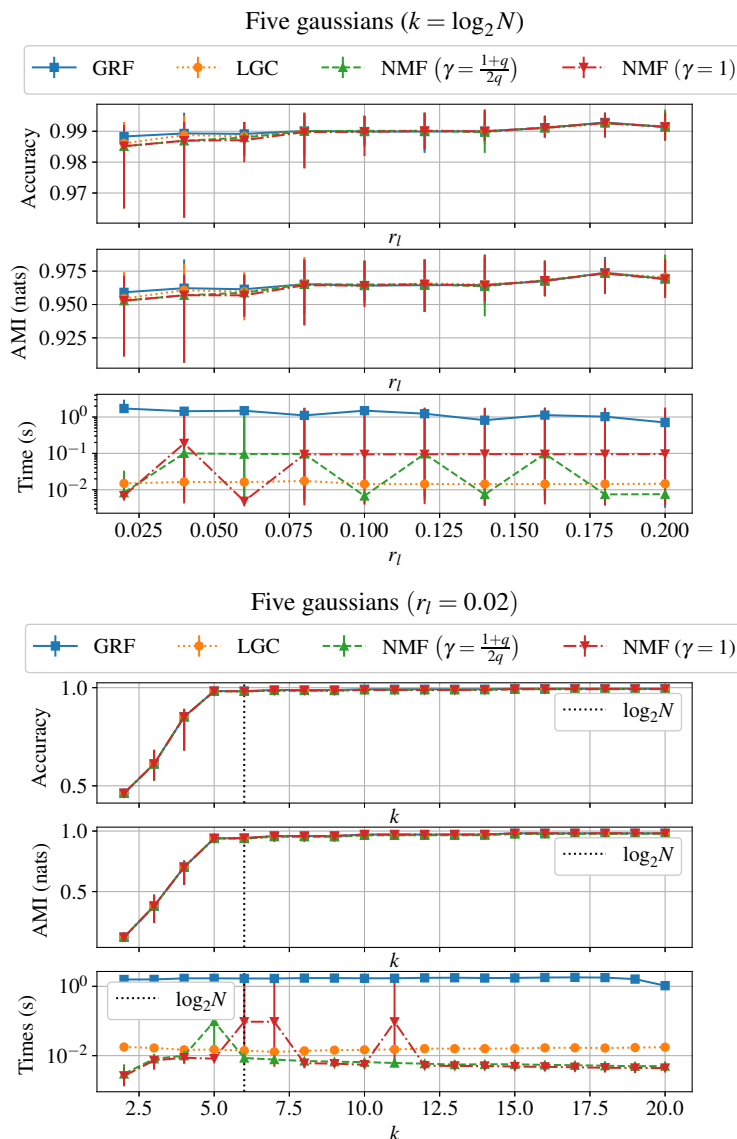


Figure 10. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

When we analyze the k dependency of the results we start observing some differences compared to the bidimensional case: a divergence between accuracy and AMI appears (particularly for $k = 5$ and $k = 6$), since $\gamma = 1$ outperforms GRF in the first metric but not in the second. As k increases past $\log_2 N$, AMI decreases for all algorithms and is followed in a less pronounced way by the accuracy.

Regarding computational performance of Digit1 (Figure 11), for $k < \log_2 N$ there is some irregularity in execution time, with the exception of NMF with $\gamma = \frac{1+q}{2q}$. With $k \geq \log_2 N$ algorithms stabilize and have similar performances.

Results on the Twonorm (Figure 12) dataset are closer to what is expected, based on the discussion on bidimensional datasets. NMF with $\gamma = 1$, LGC and GRF have similar behavior in their accuracy and AMI curves as a function of r_l . The accuracy of

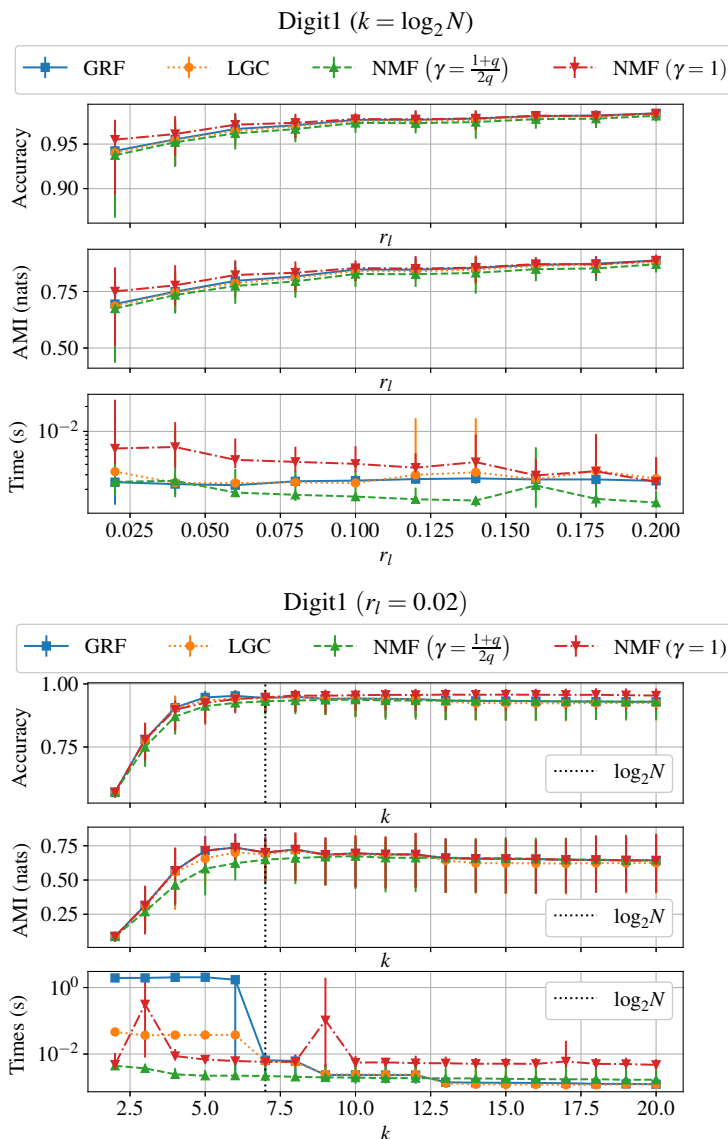


Figure 11. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

$\gamma = \frac{1+q}{2q}$ is only slightly worse than the other algorithms in absolute terms, while AMI for this method is around 0.04 smaller than other methods. However, this approach showed to be consistently faster both as a function of r_l and as a function of k . In fact, when accuracy and AMI are analyzed as functions of k , the difference between methods becomes even less relevant.

On Ringnorm (Figure 13) setting $\gamma = \frac{1+q}{2q}$ produces the most accurate results. In situations where labeled data is more scarce, this method also scores the higher values for AMI, but as r_l increases it is overcome by GRF and $\gamma = 1$. However, when we consider different constructions of G_D , both NMF approaches are more accurate than GRF and LGC by a significant margin when $k \geq \log_2 N$. In the same range, $\gamma = \frac{1+q}{2q}$ has a higher AMI than every other method and is also faster by at least one order of

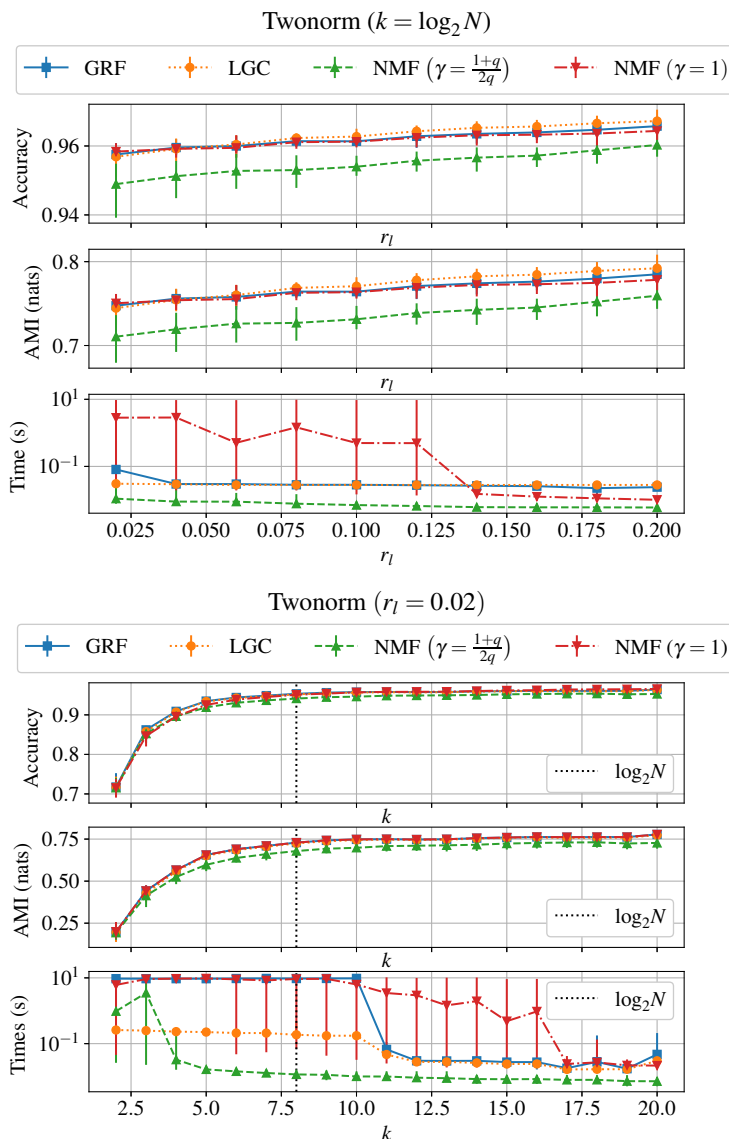


Figure 12. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

magnitude.

Landsat (Figure 14) is a more well-behaved case regarding accuracy and AMI, with all methods showing very close results. We note that LGC seems to become less susceptible to the addition of labeled data as $r_l > 0.15$. It is also noteworthy that LGC is the fastest approach in this dataset by two orders of magnitude, with exception of $k > 17$ where GRF has a significant improvement in execution time.

Experiments in the Texture dataset (Figure 15) showed some divergence among methods in the region of lower r_l , with GRF being the better method regarding accuracy and AMI, followed by NMF and then LGC. However, increasing the amount of labeled data makes the algorithms indistinguishable regarding those metrics. We also note that LGC and GRF are around two orders of magnitude faster than NMF for $k \geq \log_2 N$.

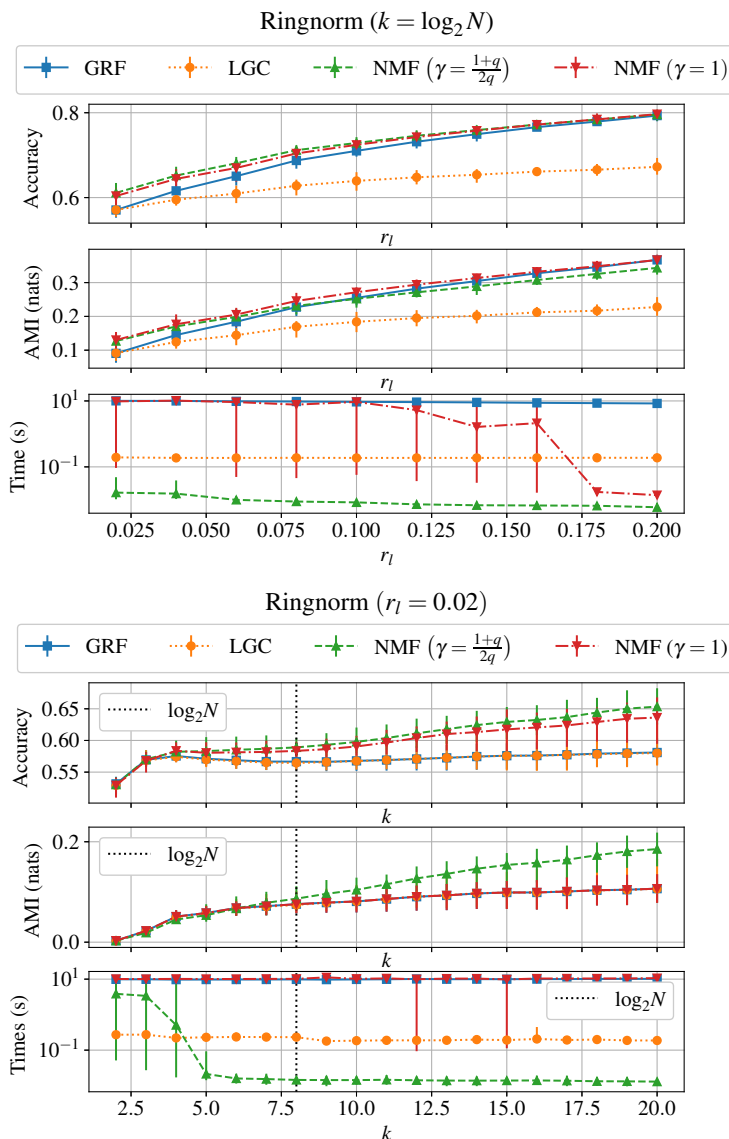


Figure 13. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

When we look at the results as a function of k for Texture, LGC shows the same behavior we highlighted for the three clusters dataset (Figure 9). For $k > \log_2 N$ accuracy and AMI suffer a significant drop, indicating that a denser graph may blur its ability to distinguish between different classes and clusters.

Finally, for USPS (Figure 16) NMF is again the slower method, followed once more by GRF and LGC. Accuracy and AMI results as a function of r_l behave similarly to the Texture dataset (Figure 15), but the divergence between methods fades at higher values of r_l .

When we analyze results in USPS as a function of k we observe a very rich behavior by all algorithms. Both NMF approaches show basically the same accuracy curve, while the case $\gamma = 1$ is much more similar to GRF in terms of AMI. In fact, these two

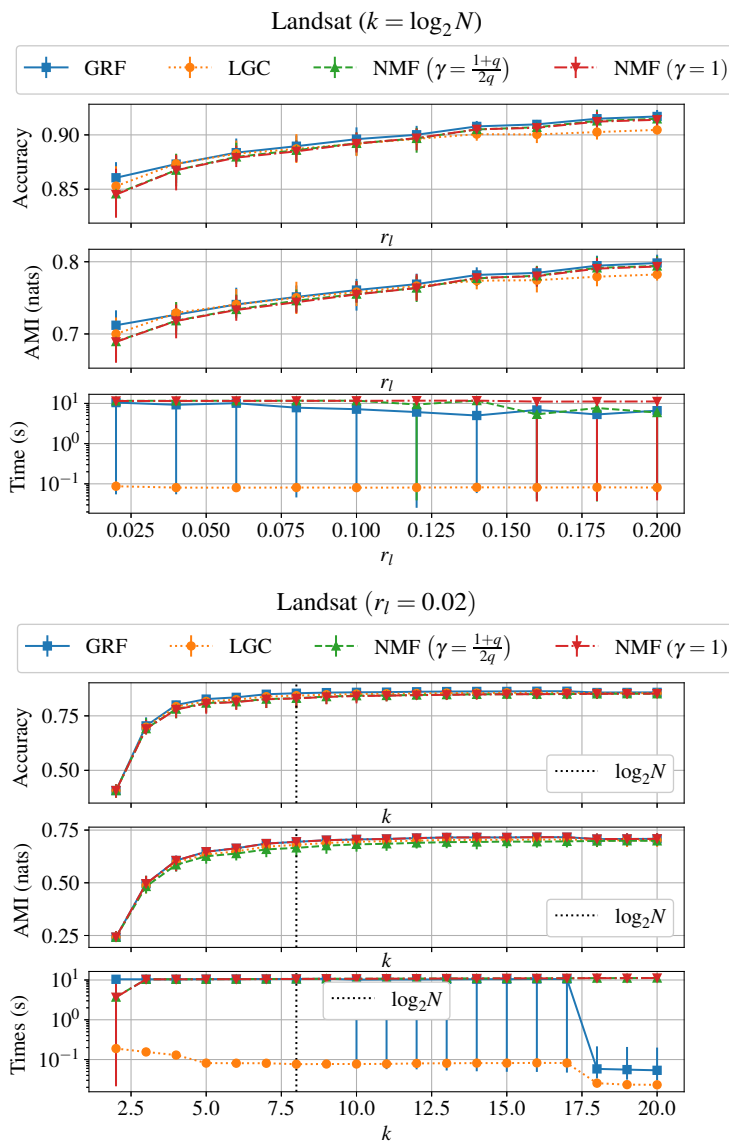


Figure 14. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

approaches show a dip in AMI for $k \geq \log_2 N$, while the accuracy of NMF, in this case, is smoother. LGC, on the other hand, shows two different falls in accuracy and LGC at different values of k , a phenomenon that also occurs in the execution time, as in GRF. The case of NMF with $\gamma = \frac{1+q}{2q}$ draws attention in this dataset due to its monotonicity with respect to k , while other methods are more susceptible to different constructions of G_D .

6.3. Discussion

Bidimensional datasets (Figures 8-10) showed a very intimate connection between accuracy and AMI, as was expected due to the clustering hypothesis of SSL. On higher

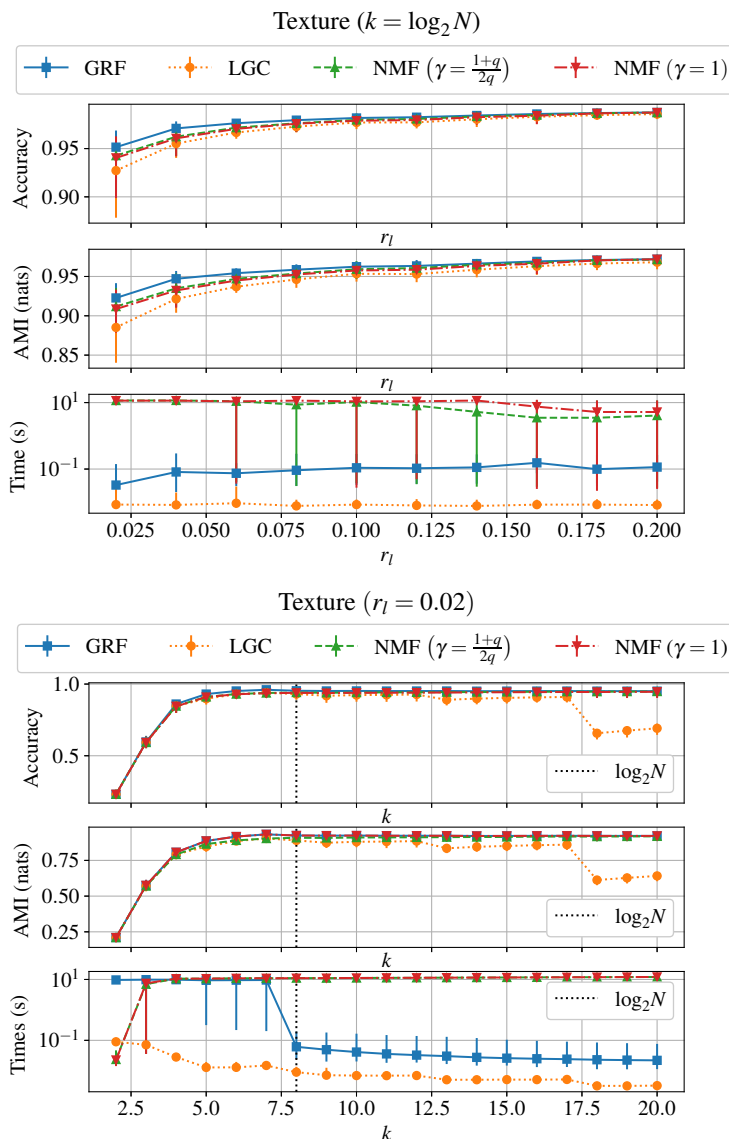


Figure 15. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

dimensional datasets (Figures 11-16), however, the connection between these metrics is not as intimate, as it was observed in Digit1 (Figure 11), Ringnorm (Figure 13), Texture (Figure 15) and USPS (Figure 16). The exact reason for this behavior is unclear to us and demands a deeper investigation.

Now, regarding our tuning method, we see that setting $\gamma = \frac{1+q}{2q}$ can lead to a faster execution time than $\gamma = 1$ as one would expect from our study in the previous section since the first approach produces lower values of β_γ^* (Figure 7). This behavior is observed in a more pronounced way in Two moons (Figure 8), Digit1 (Figure 11), Twonorm (Figure 12) and Ringnorm (Figure 13). In other datasets, both methods have a similar execution time, which leads us to believe both values of β_γ^* fall into a plateau of bad computational performance.

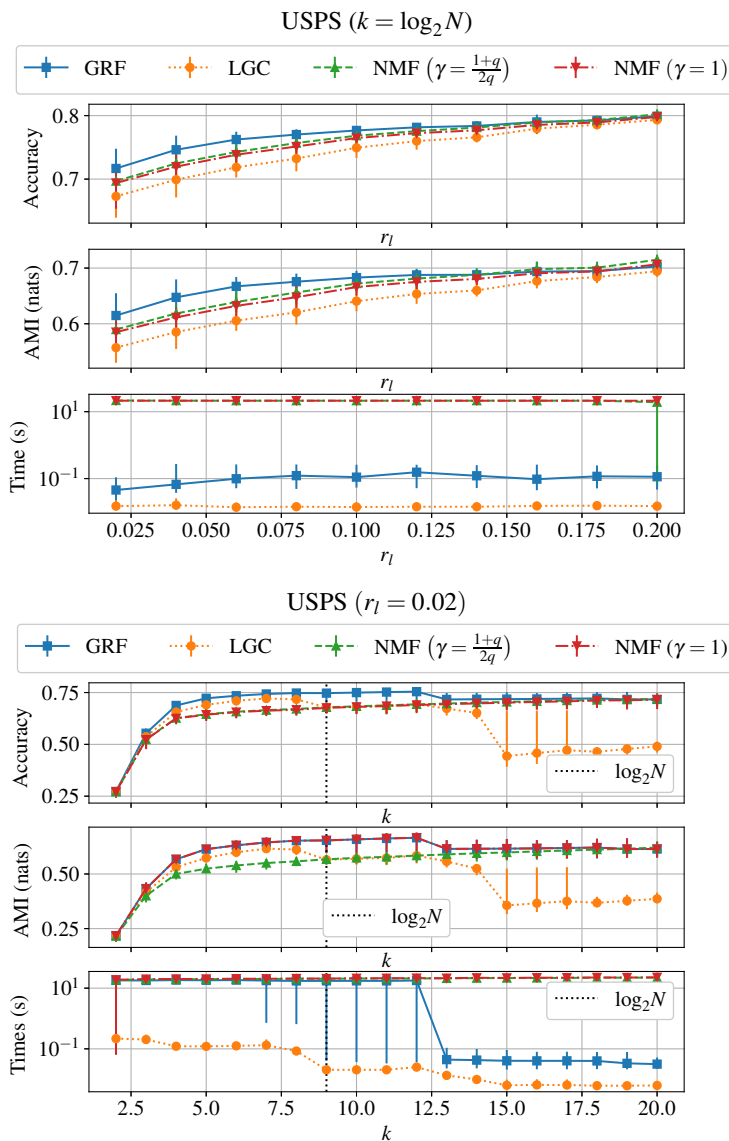


Figure 16. Results for GRF, LGC and NMF as functions of k and r_l . Lines denote averages and bars denote maximum and minimum over different realizations of D_l .

The former behavior may be related to the limitations of the approximation constructed in the previous chapter, since higher values of q lead to higher values of β_γ^* (Figure 7). In fact, for bidimensional datasets, which have $q \leq 5$, NMF shows its better performance. In the high-dimensional case, NMF with $\gamma = \frac{1+q}{2q}$ outperforms other approaches consistently when $q = 2$ (Figures 11-13). As the number of classes increases, values of β_γ^* fall in the region of worst performance (Figures 14-16) and LGC becomes the fastest algorithm. Also, results presented in section 5 suggest that this slowing down cannot be overcome by addition of labeled data, as most datasets tend to become slower with the increase in r_l .

When we analyze results on accuracy and AMI as functions of r_l we observe all algorithms have a similar behavior of improving those metrics with the addition of

labeled data. They diverge, however, on how susceptible they are to this increase, as can be seen in sets like Ringnorm (Figure 13) and Landsat (Figure 14).

Looking at the k -dependency of accuracy and AMI there is a well-defined tendency of increasing these quantities for lower k , which can then be followed by a decrease that is much less sensitive to variations in the number of nearest neighbors, as is the case of Digit1 (Figure 11). On Texture (Figure 15) and USPS (Figure 16), however, we see that different algorithms respond differently to different topologies of G_D , with Potts-based methods like GRF and NMF being more tolerant to changes in k in terms of accuracy and AMI.

It is also noteworthy that setting $\gamma = \frac{1+q}{2q}$ for NMF showed a more stable profile of accuracy and AMI as a function of k than other methods. This can be a property of interest, as choosing k for an application is a hard problem overall due to the aforementioned behavior of other algorithms.

7. Conclusion

We have studied the problem of tuning β in the NMF equations for the Potts model in applications of semi-supervised transductive classification. Through an analysis of different quantities related to the problem and the model, we were able to verify the difficulty of the problem, as optimal results are usually associated with higher computational times. As labeled data becomes scarcer, finding the best classifications can become even harder due to a well-pronounced peak in quantities like accuracy and AMI.

By the analysis of the probability of the most probable configuration, however, we were able to identify more stable classifications with higher probabilities. By using an approximation for Γ we then tested two tuning methods by using two different target values $\gamma = \frac{1+q}{2q}$ and $\gamma = 1$ for the said quantity. Results then showed that proposed approaches are effective and can improve on classical algorithms like GRF and LGC, particularly on datasets with fewer classes.

Our studies also raises questions to the possibility of achieving better computational performance in datasets with higher q while maintaining the good metrics on accuracy and AMI. This of course demands novel tuning methods to be proposed as well as the study of novel models for the task at hand.

Regarding the area of SSL, our most interesting contribution was the observation that setting $\gamma = \frac{1+q}{2q}$ leads to a more smooth dependency of accuracy and AMI as a function of k when compared to the other methods, which might be of interest to practitioners in this field.

Overall, our work helps illustrate the problem of parameter tuning in the studied model for SSL. We hope our efforts draw the attention of more researchers to this interesting problem, as this paper in no way can claim it has conquered it. We want, however, to keep our eyes on it and elaborate novel approaches as well as study other approximations and models in order to advance our understanding.

Acknowledgments

This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brasil (CAPES) – Finance Code 001. We also thank professors Denis Salvadeo from IGCE/UNESP, Alexandre Levada from DC/UFSCar and the anonymous reviewer of a previous draft of this work for insightful comments and discussions.

- [1] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. *Semi-supervised learning*. MIT Press, 2006.
- [2] Jesper E Van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine Learning*, 109(2):373–440, 2020.
- [3] Ethem Alpaydin. *Introduction to machine learning*. MIT press, 2004.
- [4] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [5] Yanwen Chong, Yun Ding, Qing Yan, and Shaoming Pan. Graph-based semi-supervised learning: A review. *Neurocomputing*, 408:216–230, 2020.
- [6] Fabricio Breve. Interactive image segmentation using label propagation through complex networks. *Expert Systems With Applications*, 123:18–33, 2019.
- [7] Peng Zhang, Tao Zhuo, Yanning Zhang, Dapeng Tao, and Jun Cheng. Online tracking based on efficient transductive learning with sample matching costs. *Neurocomputing*, 175:166–176, 2016.
- [8] Juhua Liu, Qihuang Zhong, Yuan Yuan, Hai Su, and Bo Du. Semitext: Scene text detection with semi-supervised learning. *Neurocomputing*, 407:343–353, 2020.
- [9] Yizhuo Zhang, Zhirong Wu, Houwen Peng, and Stephen Lin. A transductive approach for video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6949–6958, 2020.
- [10] Fa-Yueh Wu. The potts model. *Reviews of modern physics*, 54(1):235, 1982.
- [11] Thomas Ott, Albert Kern, Ausgar Schuffenhauer, Maxim Popov, Pierre Acklin, Edgar Jacoby, and Ruedi Stoop. Sequential superparamagnetic clustering for unbiased classification of high-dimensional chemical data. *Journal of chemical information and computer sciences*, 44(4):1358–1364, 2004.
- [12] Fei Wang, Shijun Wang, Changshui Zhang, and Ole Winther. Semi-supervised mean fields. In *Artificial Intelligence and Statistics*, pages 596–603, 2007.
- [13] Jianqiang Li and Fei Wang. Semi-supervised learning via mean field methods. *Neurocomputing*, 177:385–393, 2016.
- [14] Marcelo Blatt, Shai Wiseman, and Eytan Domany. Superparamagnetic clustering of data. *Physical review letters*, 76(18):3251, 1996.
- [15] Gad Getz, Noam Shental, and Eytan Domany. Semi-supervised learning—a statistical physics approach. In *“Learning with Partially Classified Training Data”—ICML05 workshop*. Citeseer, 2005.
- [16] Masayuki Karasuyama and Hiroshi Mamitsuka. Adaptive edge weighting for graph-based learning algorithms. *Machine Learning*, 106(2):307–335, 2017.
- [17] Junliang Ma, Bing Xiao, and Cheng Deng. Graph based semi-supervised classification with probabilistic nearest neighbors. *Pattern Recognition Letters*, 133:94–101, 2020.
- [18] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label propagation. 2002.
- [19] Xiaojin Zhu and John Lafferty. Harmonic mixtures: combining mixture models and graph-based methods for inductive and scalable semi-supervised learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 1052–1059, 2005.
- [20] Dengyong Zhou, Olivier Bousquet, Thomas N Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. In *Advances in neural information processing*

- systems*, pages 321–328, 2004.
- [21] Tony Jebara, Jun Wang, and Shih-Fu Chang. Graph construction and b-matching for semi-supervised learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 441–448, 2009.
 - [22] Celso André R de Sousa, Solange O Rezende, and Gustavo EAPA Batista. Influence of graph construction on semi-supervised learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 160–175. Springer, 2013.
 - [23] Amarnag Subramanya and Partha Pratim Talukdar. Graph-based semi-supervised learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 8(4):1–125, 2014.
 - [24] Gergely Tibély and János Kertész. On the equivalence of the label propagation method of community detection and a potts model approach. *Physica A: Statistical Mechanics and its Applications*, 387(19-20):4982–4984, 2008.
 - [25] Cheng Shi, Yanchen Liu, and Pan Zhang. Weighted community detection and data clustering using message passing. *Journal of Statistical Mechanics: Theory and Experiment*, 2018(3):033405, 2018.
 - [26] Manfred Opper and David Saad. *Advanced mean field methods: Theory and practice*. MIT press, 2001.
 - [27] Delve: Data for evaluating learning in valid experiments. 1995.
 - [28] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
 - [29] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
 - [30] Jesús Alcalá-Fdez, Alberto Fernández, Julián Luengo, Joaquín Derrac, Salvador García, Luciano Sánchez, and Francisco Herrera. Keel data-mining software tool: data set repository, integration of algorithms and experimental analysis framework. *Journal of Multiple-Valued Logic & Soft Computing*, 17, 2011.
 - [31] Nguyen Xuan Vinh, Julien Epps, and James Bailey. Information theoretic measures for clusterings comparison: Variants, properties, normalization and correction for chance. *The Journal of Machine Learning Research*, 11:2837–2854, 2010.
 - [32] Dennis M Ritchie, Brian W Kernighan, and Michael E Lesk. *The C programming language*. Prentice Hall Englewood Cliffs, 1988.
 - [33] Stefan Behnel, Robert Bradshaw, Craig Citro, Lisandro Dalcin, Dag Sverre Seljebotn, and Kurt Smith. Cython: The best of both worlds. *Computing in Science & Engineering*, 13(2):31–39, 2010.
 - [34] Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009.
 - [35] Leonardo Dagum and Ramesh Menon. Openmp: an industry standard api for shared-memory programming. *IEEE computational science and engineering*, 5(1):46–55, 1998.
 - [36] Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods*, 17(3):261–272, 2020.