
On Pairwise Clustering with Side Information

Stephen Pasteris

Department of Computer Science
University College London
London, UK
S.Pasteris@cs.ucl.ac.uk

Fabio Vitale

INRIA Lille
Lille, France
fabio.vitale@inria.fr

Claudio Gentile

DiSTA
University of Insubria
Varese, Italy
claudio.gentile@uninsubria.it

Mark Herbster

Department of Computer Science
University College London
London, UK
M.Herbster@cs.ucl.ac.uk

Abstract

Pairwise clustering, in general, partitions a set of items via a known similarity function. In our treatment, clustering is modeled as a transductive prediction problem. Thus rather than beginning with a known similarity function, the function instead is hidden and the learner only receives a random sample consisting of a subset of the pairwise similarities. An additional set of pairwise side-information may be given to the learner, which then determines the inductive bias of our algorithms. We measure performance not based on the recovery of the hidden similarity function, but instead on how well we classify each item. We give tight bounds on the number of misclassifications. We provide two algorithms. The first algorithm *SACA* is a simple agglomerative clustering algorithm which runs in near linear time, and which serves as a baseline for our analyses. Whereas the second algorithm, *RGCA*, enables the incorporation of side-information which may lead to improved bounds at the cost of a longer running time.

1 Introduction

The aim of clustering is to partition a set of n items into k “clusters” based on their similarity. A common approach to clustering is to assume that items can be embedded in a metric space, and then to (approximately) minimize an objective function over all possible partitionings based on the metric at hand. A quintessential example is the k -means objective. An alternative is to assume only the existence of a similarity function between the pairs. Examples of this approach include spectral [19] and k -median [18], as well as correlation clustering [2].

Our approach to clustering is most similar to correlation clustering. Correlation clustering was introduced in the seminal paper [2]. In this setting, a complete graph of similarity and dissimilarity item pairs is given. The goal is to find a disjoint partition (“clustering”) which minimizes an objective that counts the total number of “incorrect” similarity and dissimilarity pairs in the resulting clustering. A pair of items is incorrect with respect to the clustering if it was given as similar while they appear in distinct clusters, and vice versa. In [2] an efficient algorithm with a guaranteed approximation ratio was given for this NP-hard problem.

Although inspired by these results, our focus is slightly different: we seek to provide efficient algorithms that compute a clustering, as well as to provide predictive performance guarantees for these algorithms.

We treat pairwise clustering as a transductive prediction problem. Given a set of unlabeled items, the aim is to predict their class labels. As input to our algorithms, firstly we have a training set of similarity and dissimilarity item pairs. Secondly, we have a set of *soft* similarity pairwise constraints – the *side-information* graph. The side-information graph determines our inductive bias, i.e., our output clustering will tend to (but need not) place each softly constrained item pair into the same cluster. We give bounds for a batch learning model where the learner samples uniformly at random a training set of m similarity/dissimilarity pairs from a ground truth clustering. Given these m pairs and the side-information graph, the learner then outputs a clustering. The quality of the resulting clustering is measured by the item *misclassification error* which is essentially the number of items in the learner’s output clustering that are misclassified as compared to the ground truth. We describe and analyze two novel algorithms for pairwise clustering, and deliver upper bounds on their expected misclassification error which scale to the degree that a clustering exists that reflects the inductive bias induced by the side-information graph at hand. We complement our upper bound with an almost matching lower bound on the prediction complexity of this problem.

The paper is organized as follows. In Section 2 we review notation as well as formally introduce our learning models. In Section 3, we present our two clustering algorithms RGCA and SACA, along with their analyses. Our fastest algorithm is quite efficient, for it requires only a linear time in the input size up to a sub-logarithmic factor to compute a clustering with small error. We give concluding remarks in Section 4. Finally, below we provide pointers to a few references in distinct but closely related research areas.

Related work

The literature most directly related to our work in perspective is the literature on clustering with side information, as well the literature on semi-supervised clustering. Some of the references in this area include [3, 25]. Secondly, our work is also connected to the metric learning task. Metric learning is also concerned with recovering a similarity function; however, in this literature the similarity is treated as a real-valued function often identified with a positive semi-definite matrix as opposed to our binary model. Some relevant references here include [30, 21, 4]. What distinguishes our work from the past literature is that we are aimed at constructing clusterings with side information, not just similarity functions, with an associated tight misclassification error analysis.

2 Preliminaries and Notation

We now introduce our main notation along with basic preliminaries. Given a finite set $V = \{1, \dots, n\}$, a *clustering* \mathcal{D} over V is a partition of V into a finite number of sets $\mathcal{D} = \{D_1, \dots, D_k\}$. Each D_j is called a *cluster*. A *similarity* graph $G = (V, \mathcal{P})$ over V is an undirected (but not necessarily connected) graph where, for each pairing $(v, w) \in V^2$, v and w are *similar* if $(v, w) \in \mathcal{P}$, and *dissimilar*, otherwise. Notice that the similarity relationship so defined need not be transitive. We shall interchangeably represent a similarity graph over V through a binary $n \times n$ *similarity* matrix $Y = [y_{v,w}]_{v,w=1}^{n \times n}$ whose entry $y_{v,w}$ is 1 if items v and w are similar, and $y_{v,w} = 0$, otherwise. A clustering \mathcal{D} over V can be naturally associated with a similarity graph $G = (V, \mathcal{P}_{\mathcal{D}})$ whose edge set $\mathcal{P}_{\mathcal{D}}$ is defined as follows: Given $v, w \in V$, then $(v, w) \in \mathcal{P}_{\mathcal{D}}$ if and only if there exists a cluster $D \in \mathcal{D}$ with $v, w \in D$. In words, G is made up of k disjoint cliques. It is only in this case that the similarity relationship defined through G is transitive. Matrix Y represents a clustering if, after permutation of rows and columns, it ends up being block-diagonal, where the i -th block is a $d_i \times d_i$ matrix of ones, d_i being the size of the i -th cluster. Given clustering \mathcal{D} , we find it convenient to define a map $\mu_{\mathcal{D}} : V \rightarrow \{1, \dots, k\}$ in such a way that for all $v \in V$ we have $v \in D_{\mu_{\mathcal{D}}(v)}$. In words, $\mu_{\mathcal{D}}$ is a class assignment mapping, so that v and w are similar w.r.t. \mathcal{D} if and only if $\mu_{\mathcal{D}}(v) = \mu_{\mathcal{D}}(w)$.

Given two similarity graphs $G = (V, \mathcal{P})$ and $G' = (V, \mathcal{P}')$, the (Hamming error) distance between G and G' , denoted here as $\text{HA}(\mathcal{P}, \mathcal{P}')$, is defined as

$$\text{HA}(\mathcal{P}, \mathcal{P}') = \left| \{(v, w) \in V^2 : (v, w) \in \mathcal{P} \wedge (v, w) \notin \mathcal{P}' \vee (v, w) \notin \mathcal{P} \wedge (v, w) \in \mathcal{P}'\} \right|,$$

where $|A|$ is the cardinality of set A . The same definition applies in particular to the case when either G or G' (or both) represent clusterings over V . By abuse of notation, if \mathcal{D} is a clustering and $G = (V, \mathcal{P})$ is a similarity graph, we will often write $\text{HA}(\mathcal{D}, \mathcal{P})$ to denote $\text{HA}(\mathcal{P}_{\mathcal{D}}, \mathcal{P})$, where $(V, \mathcal{P}_{\mathcal{D}})$ is the similarity graph associated with \mathcal{D} , so that $\text{HA}(\mathcal{P}_{\mathcal{D}}, \mathcal{D}) = 0$. Moreover, if the similarity graphs G

and G' are represented by similarity matrices, we may equivalently write $\text{HA}(Y, Y')$, $\text{HA}(Y, \mathcal{D})$, and so on. The quantity $\text{HA}(\cdot, \cdot)$ is very closely related to the so-called Mirkin metric [24] over clusterings, as well as to the (complement of the) Rand index [27], see, e.g., [23].

Another “distance” that applies specifically to clusterings is the misclassification error distance, denoted here as $\text{ER}(\cdot, \cdot)$, and is defined as follows. Given two clusterings $\mathcal{C} = \{C_1, \dots, C_\ell\}$ and $\mathcal{D} = \{D_1, \dots, D_k\}$ over V , repeatedly add the empty set to the smaller of the two so as to obtain $\ell = k$. Then

$$\text{ER}(\mathcal{C}, \mathcal{D}) = \min_f \sum_{D \in \mathcal{D}} |D \setminus f(D)|,$$

the minimum being over all bijections from \mathcal{D} to \mathcal{C} . In words, $\text{ER}(\mathcal{C}, \mathcal{D})$ measures the smallest number of classification mistakes over all class assignments of clusters in \mathcal{D} w.r.t. clusters in \mathcal{C} . This is basically an unnormalized version of the classification error distance considered, e.g., in [22].

The (Jaccard) distance $\text{DIST}(A, B)$ between sets A and B , with $A, B \subseteq V$ is defined as

$$\text{DIST}(A, B) = \frac{|A \setminus B| + |B \setminus A|}{|A \cup B|}.$$

Recall that $\text{DIST}(\cdot, \cdot)$ is a proper metric on the collection of all finite sets. Moreover, observe that $\text{DIST}(A, B) = 1$ if and only if A and B are disjoint.

Since our clustering algorithms will rely upon side information in the form of undirected graphs, we also need to recall relevant notions for such graphs and (spectral) properties thereof. Let Y be a similarity matrix and $G = (V, E)$ be a graph, henceforth called *side-information* graph. G is assumed to be undirected, unweighted and connected.

As is standard in graph-based learning problems (e.g., [12–14, 16, 15, 6–8, 28, 10], and references therein), graph G encodes side information in that it suggests to the clustering algorithms that adjacent vertices in G tend to be similar. The set of *cut-edges* in G w.r.t. Y is the set of edges $(v, w) \in E$ such that $y_{v,w} = 0$, the associated cut-size (i.e., their number) will be denoted as $\Phi_G(Y)$ (or simply Φ_G , if Y is clear from the surrounding context).

If G is viewed as a resistive network where each edge is a unit resistor, then the *effective resistance* $r_G(v, w)$ of the pairing $(v, w) \in V^2$ is a measure of connectivity between the two nodes v and w in G which, in the special case when $(v, w) \in E$, also equals the probability that a spanning tree of G drawn uniformly at random from the set of all spanning trees of G includes (v, w) as one of its $n - 1$ edges (e.g., [20]). As a consequence, $\sum_{(v,w) \in E} r_G(v, w) = n - 1$. Finally, $\Phi_G^R(Y)$ (or Φ_G^R , for brevity) will denote the sum, over all cut-edges (v, w) in G w.r.t. Y , of the effective resistances $r_G(v, w)$. This sum will sometimes be called the *resistance-weighted* cut-size of G (w.r.t. Y). Notice that if G is a tree we have $\Phi_G^R(Y) = \Phi_G(Y)$ for all Y .

The basic inductive principle underpinning RGCA is the assumption that $\Phi_G^R(Y)$ is small.¹ Both Φ_G^R and Φ_G can be considered as complexity measures for our learning problems, since they both depend on cut-edges in E . However, unlike Φ_G , the quantity Φ_G^R enjoys properties of *global* density-independence (Φ_G^R is at most $n - 1$, hence it scales with the number of nodes of G rather than the number of edges), and *local* density-independence (Φ_G^R suitably discriminates between dense and sparse graph topology areas – see, e.g., the discussion in [7]). As such, Φ_G^R is more satisfactory than Φ_G in measuring the quality of side information at our disposal.

2.1 Learning setting

We are interested in inferring (or just computing) clusterings over V based on binary similarity/dissimilarity information contained in a similarity matrix Y , possibly along with side information in the form of a connected and undirected graph $G = (V, E)$. The similarity matrix Y itself may or may not represent a clustering over V . The error of our inference procedures will be measured through $\text{ER}(\cdot, \cdot)$. We shall find bounds on $\text{ER}(\cdot, \cdot)$ either directly, by presenting specific algorithms, or indirectly via (tight) reductions from similarity prediction problems/methods measured through $\text{HA}(\cdot, \cdot)$ to clustering problems/methods measured through $\text{ER}(\cdot, \cdot)$. More specifically, given a set of

¹ Notice that the edges in G should not be considered as hard constraints (like the must-link constraints in semi-supervised clustering/clustering with side information, e.g., [3, 9]).

Algorithm 1 The Robust Greedy Clustering Algorithm

Input: Similarity graph (V, \mathcal{P}) ; distance parameter $a \in [0, 1]$.

1. For all $v \in V$, set $\Gamma(v) \leftarrow \{v\} \cup \{w \in V : (v, w) \in \mathcal{P}\}$;
2. Construction of graph (V, \mathcal{Q}) : //First stage
For all $v, w \in V$ with $v \neq w$:
If $\text{DIST}(\Gamma(v), \Gamma(w)) \leq 1 - a$ then $(i, j) \in \mathcal{Q}$, otherwise $(i, j) \notin \mathcal{Q}$;
3. Set $A_1 \leftarrow V$, and $t \leftarrow 1$; //Second stage
4. While $A_t \neq \emptyset$:
 - For every $v \in A_t$ set $N_t(v) \leftarrow \{v\} \cup \{w \in A_t : (v, w) \in \mathcal{Q}\}$,
 - Set $\alpha_t \leftarrow \operatorname{argmax}_{v \in A_t} |N_t(v)|$,
 - Set $C_t \leftarrow N_t(\alpha_t)$,
 - Set $A_{t+1} \leftarrow A_t \setminus C_t$,
 - $t \leftarrow t + 1$;

Output: C_1, C_2, \dots, C_ℓ , where $\ell = t - 1$.

items $V = \{1, \dots, n\}$ and a similarity matrix Y representing a clustering \mathcal{D} , our goal is to build a clustering \mathcal{C} over V with as small as possible $\text{ER}(\mathcal{C}, \mathcal{D})$. We would like to do so by observing only a subset of the binary entries of Y . Notice that the number of clusters k in the comparison clustering need not be known to the clustering algorithm.

In the setting of RGCA, we are given a side information graph $G = (V, E)$, and a training set S of m binary-labeled pairs $\langle (v, w), y_{u,v} \rangle \in V^2 \times \{0, 1\}$, drawn uniformly at random² from V^2 . Our goal is to build a clustering \mathcal{C} over V so as to achieve small misclassification error $\text{ER}(\mathcal{C}, Y)$, when this error is computed *on the whole* matrix Y .

3 Algorithms and Analysis

We start off with a clustering algorithm that takes as input a similarity graph over V , and produces in output a clustering over V . This will be a building block for later results, but it can also be of independent interest. Our algorithm, called Robust Greedy Clustering Algorithm (RGCA, for brevity), is displayed in Algorithm 1. The algorithm has two stages. The first stage is a robustifying stage where the similarity graph (V, \mathcal{P}) is converted into a (more robust) similarity graph (V, \mathcal{Q}) as follows: Given two distinct vertices $v, w \in V$, we have $(v, w) \in \mathcal{Q}$ if and only if the Jaccard distance of their neighbourhoods (in (V, \mathcal{P})) is not bigger than $1 - a$, for some distance parameter $a \in [0, 1]$. The second stage uses a greedy method to convert the graph (V, \mathcal{Q}) into a clustering \mathcal{C} . This stage proceeds in “rounds”. At each round t we have a set A_t of all vertices which have not yet been assigned to any clusters. We then choose α_t to be the vertex in A_t which has the maximum number of neighbours (under the graph (V, \mathcal{Q})) in A_t , and take this set of neighbours (including α_t) to be the next cluster.

From a computational standpoint, the second stage of RGCA runs in $\mathcal{O}(n^2 \log n)$ time, since on every round t we single out α_t (which can be determined in $\log n$ time by maintaining a suitable heap data-structure), and erase all edges emanating from α_t in the similarity graph (V, \mathcal{Q}) . On the other hand, the first stage of RGCA runs in $\mathcal{O}(n^3)$ time, in the worst case, though standard techniques exist that avoid the all-pairs comparison, like a Locality Sensitive Hashing scheme applied to the Jaccard distance (e.g., [26, Ch.3]). We have the following result.³

Theorem 1. *Let $\mathcal{C} = \{C_1, \dots, C_k\}$ be the clustering produced in output by RGCA when receiving as input similarity graph (V, \mathcal{P}) , and distance parameter $a = 2/3$. Then for any clustering $\mathcal{D} = \{D_1, \dots, D_k\}$, with $d_i = |D_i|$, $i = 1, \dots, k$, and $d_1 \leq d_2 \leq \dots \leq d_k$ we have*

$$\text{ER}(\mathcal{C}, \mathcal{D}) \leq \min_{j=1, \dots, k} \left(\frac{12}{d_j} \text{HA}(\mathcal{P}, \mathcal{D}) + \sum_{i=1}^{j-1} d_i \right).$$

² For simplicity of presentation, we will assume the samples in S are drawn from V with replacement.

³ All proofs are contained in the appendix.

Hence, if the chosen \mathcal{D} is the best approximation to \mathcal{P} w.r.t. $\text{HA}(\cdot, \cdot)$, and we interpret (V, \mathcal{P}) as a noisy version of \mathcal{D} , then small $\text{HA}(\mathcal{P}, \mathcal{D})$ implies small $\text{ER}(\mathcal{C}, \mathcal{D})$. In particular, $\text{HA}(\mathcal{P}, \mathcal{D}) = 0$ implies $\text{ER}(\mathcal{C}, \mathcal{D}) = 0$ (simply pick $j = 1$ in the minimum). Yet, this result only applies to the case when the similarity graph (V, \mathcal{P}) is fully observed by our clustering algorithm. As we will see below, (V, \mathcal{P}) may in turn be the result of a similarity learning process when the similarity labels are provided by clustering \mathcal{D} . In this sense, Theorem 1 will help us to deliver generalization bounds (as measured by $\text{ER}(\mathcal{C}, \mathcal{D})$), as a function of the generalization ability of this similarity learning process (as measured by $\text{HA}(\mathcal{P}, \mathcal{D})$).

The problem faced by RGCA is also related to the standard correlation clustering problem [2]. Yet, the goal here is somewhat different, since a correlation clustering algorithm takes as input (V, \mathcal{P}) , but is aimed at producing a clustering \mathcal{C} such that $\text{HA}(\mathcal{P}, \mathcal{C})$ is as small as possible.

In passing, we next show that the construction provided by RGCA is essentially optimal (up to multiplicative constants). Let $G_{\mathcal{D}} = (V, E_{\mathcal{D}})$ be the similarity graph associated with clustering \mathcal{D} . We say that a clustering algorithm that takes as input a similarity graph over V and gives in output a clustering over V is *consistent* if and only if for every clustering \mathcal{D} over V the algorithm outputs \mathcal{D} when receiving as input $G_{\mathcal{D}}$. Observe that RGCA is an example of a consistent algorithm. We have the following lower bound.

Theorem 2. *For any finite set V , any clustering $\mathcal{D} = \{D_1, D_2, \dots, D_k\}$ over V , any positive constant σ , and any consistent clustering algorithm, there exists a similarity graph (V, \mathcal{P}) such that $\text{HA}(\mathcal{P}, \mathcal{D}) \leq \sigma$, while*

$$\text{ER}(\mathcal{C}, \mathcal{D}) \geq \min_{j=1, \dots, k} \left(\frac{1}{2d_j} \sigma - 1 + \frac{1}{4} \sum_{i=1}^{j-1} d_i \right), \quad (1)$$

or $\text{ER}(\mathcal{C}, \mathcal{D}) \geq \frac{n}{2}$, where \mathcal{C} is the output produced by the algorithm when given (V, \mathcal{P}) as input, and $d_i = |D_i|$, $i = 1, \dots, k$, with $d_1 \leq d_2 \leq \dots \leq d_k$.

From the proof provided in the appendix, one can see that the similarity graph (V, \mathcal{P}) used here is indeed a *clustering* over V so that, as the algorithm is consistent, the output \mathcal{C} must be such a clustering. This result can therefore be contrasted to the results contained, e.g., in [23] about the equivalence between clustering distances, specifically Theorem 26 therein. Translated into our notation, that result reads as follows: $\text{ER}(\mathcal{C}, \mathcal{D}) \geq \frac{\text{HA}(\mathcal{P}, \mathcal{D})}{16d_k}$. Our Theorem 2 is thus sharper but, unlike the one in [23], it *does not apply* to any possible pairs of clusterings \mathcal{C} and \mathcal{D} , for in our case \mathcal{C} is selected as a function of \mathcal{D} .

3.1 Learning to Cluster

Suppose now that our clustering algorithm has at its disposal a side information graph G , and a training set S of size m . Training set S is drawn at random from V^2 , and is labeled according to a similarity matrix Y representing a clustering $\mathcal{D} = \{D_1, \dots, D_k\}$ with cluster sizes $d_i = |D_i|$, $i = 1, \dots, k$, and having resistance-weighted cutsize $\Phi_G^R(Y)$. A Laplacian-regularized Matrix Winnow algorithm [29], as presented in [10], is an online algorithm that sweeps over S only once, and is guaranteed to make $\mathcal{O}(\Phi_G^R \log^3 n)$ many mistakes in expectation (see Theorem 5 therein). In turn, this algorithm can be used within an online-to-batch conversion wrapper, like the one mentioned in [11], or the one in [5] to produce a similarity graph (V, \mathcal{P}) (which need not be a clustering) such that

$$\mathbb{E} \text{HA}(\mathcal{P}, Y) = \mathcal{O} \left(\frac{n^2}{m} \Phi_G^R \log^3 n \right).$$

Algorithm 2 The Simple Agglomerative Clustering Algorithm.

Input: Item set $V = \{1, \dots, n\}$; training set S .

1. Initialization: $\mathcal{C} = \{\{1\}, \dots, \{n\}\}$;
2. For any $v \in V$, let C_v denote the cluster of \mathcal{C} containing v ;
3. For each $(v, w) \in S$:
 - If $(y_{v,w} = 1) \wedge (C_v \neq C_w)$ then $\mathcal{C} \leftarrow \mathcal{C} \setminus C_w$ and $C_v \leftarrow C_v \cup C_w$;

Output: Clustering \mathcal{C} .

Then, in order to produce a "good" clustering \mathcal{C} out of \mathcal{P} , we can apply RGCA to input (V, \mathcal{P}) . Invoking Theorem 1, we conclude that

$$\begin{aligned}
 \mathbb{E} \text{ER}(\mathcal{C}, \mathcal{D}) &\leq \mathbb{E} \left[\min_{j=1, \dots, k} \left(\frac{12}{d_j} \text{HA}(\mathcal{P}, \mathcal{D}) + \sum_{i=1}^{j-1} d_i \right) \right] \\
 &\leq \min_{j=1, \dots, k} \left(\frac{12}{d_j} \mathbb{E} \text{HA}(\mathcal{P}, \mathcal{D}) + \sum_{i=1}^{j-1} d_i \right) \\
 &= \mathcal{O} \left(\min_{j=1, \dots, k} \left(\frac{1}{d_j} \frac{n^2}{m} \Phi_G^R \log^3 n + \sum_{i=1}^{j-1} d_i \right) \right). \tag{2}
 \end{aligned}$$

The training time of the whole procedure is dominated by the $\mathcal{O}(n^3)$ time per round required by Matrix Winnow, which is thus $\mathcal{O}(mn^3)$. In what follows, we take a more direct (and time-efficient) route to obtain alternative statistical guarantees in the simplest case when the side-information graph is absent.

Algorithm 2 displays the pseudocode of SACA (Simple Agglomerative Clustering Algorithm). SACA takes as input the item set V and a training set S . The algorithm operates as follows. It starts by assigning a different cluster to each vertex in V , and sequentially inspects each $\langle (v, w), y_{v,w} \rangle \in S$ aiming to merge clusters. In particular, whenever v and w currently fall into different clusters but $y_{v,w} = 1$, the two clusters are merged, as in a standard agglomerative clustering procedure. Finally, SACA outputs the clustering \mathcal{C} so computed. Notice that, unlike the Matrix Winnow-based algorithm, no side-information in the form of a graph over V is exploited.

The following theorem quantifies the performance of SACA.

Theorem 3. *Given similarity matrix Y encoding a clustering over V with k clusters, SACA returns a clustering \mathcal{C} such that $\text{ER}(\mathcal{C}, Y)$ is bounded as*

$$\mathbb{E} [\text{ER}(\mathcal{C}, Y)] = \mathcal{O} \left(\frac{n^2}{m} k \log \frac{n^2}{m} \right),$$

the expectation being over a random draw of S .

It is instructive to compare the upper bounds contained in Theorem 3 to the one in Eq. (2). The two bounds are in general incomparable. While Φ_G^R is always at least as large as $k - 1$ (recall that G is connected), the bound in Theorem 3 does also depend in a detailed way on the sizes d_i of the underlying clustering \mathcal{D} . For instance, if $d_i = n/k$ for all i then (2) is sharper in the presence of informative side-information than the bound in Theorem 3. This is because, up to log factors, the resulting bound is of order $\frac{n}{m} \Phi_G^R$ which is no larger than the bound in Theorem 3 since $k - 1 \leq \Phi_G^R \leq n - 1$. Thus in the case of maximally informative side-information ($\Phi_G^R = \Theta(k)$) and balanced cluster sizes the bound is improved by a factor of $\frac{k}{n}$. On the other hand, SACA is definitely much faster than the Matrix Winnow-based algorithm since, apart from the random spanning tree construction, it only takes $\mathcal{O}((n + m) \log^* n)$ time to run if implemented via standard (union-find) data-structures, where $\log^* n$ is the iterated logarithm of n .

We complement the two upper bounds with the following lower bound result, showing that the dependence of $\text{ER}(\cdot)$ on Φ_G^R (or k) cannot be eliminated.

Theorem 4. Given any side-information graph $G = (V, E)$, any $b \in [4, n-1]$, any $k > 2$ and any $m < \frac{n^2}{4}$, there exists a similarity matrix Y representing a clustering formed by at most k clusters such that for any algorithm giving in output clustering \mathcal{C} we have $\text{ER}(\mathcal{C}, Y) = \Omega\left(\min\left\{\frac{n^2}{m}k, b\right\}\right)$ while $\Phi_G^R(Y) \leq b$.

4 Conclusions and Ongoing Research

We have investigated the problem of learning a clustering over a finite set from pairwise training data and side-information data. Two routes have been followed to tackle this problem: i. a direct route, where we exhibited a specific algorithm, called SACA, operating without side information graph, and ii. an indirect route that steps through a reduction, called RGCA, establishing a tight bridge between two clustering metrics, that takes the side-information graph into account. We provided two misclassification error analyses in the case when the source of similarity data is consistent with a given clustering, and complemented these analyses with an involved construction delivering an almost matching lower bound.

Two extensions we are currently exploring are: i. extending the underlying statistical assumptions on data (e.g., sampling distribution-free guarantees) while retaining running time efficiency, and ii. studying other learning regimes, like active learning, under similar or broader statistical assumptions as those currently in this paper.

Acknowledgements. This work was supported in part by the U.S. Army Research Laboratory and the U.K. Defence Science and Technology Laboratory and was accomplished under Agreement Number W911NF-16-3-0001. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Army Research Laboratory, the U.S. Government, the U.K. Defence Science and Technology Laboratory or the U.K. Government. The U.S. and U.K. Governments are authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

References

- [1] N. Alon, C. Avin, M. Koucky, G. Kozma, Z. Lotker, and M. Tuttle. Many random walks are faster than one. *C Comb. Probab. Comput.*, 20(4), pp. 481–502, 2011.
- [2] N. Bansal, A. Blum S. Chawla. Correlation Clustering. *Machine Learning*, 56/1, pp. 89–113, 2004.
- [3] A. Ben-Dor, R. Shamir, and Z. Yakhini. Clustering gene expression patterns. *Journal of Computational Biology*, 6(3/4), 1999.
- [4] Q. Cao, Z. Guo, and Y. Ying. Generalization bounds for metric and similarity learning. *CoRR*, abs/1207.5437, 2012.
- [5] N. Cesa-Bianchi, C. Gentile. Improved risk tail bounds for on-line algorithms. *Information Theory, IEEE Transactions on*, 54/1, pp. 386–390.
- [6] N. Cesa-Bianchi, C. Gentile, and F. Vitale. Fast and optimal prediction on a labeled tree. In *Proceedings of the 22nd Annual Conference on Learning*. Omnipress, 2009.
- [7] N. Cesa-Bianchi, C. Gentile, F. Vitale, and G. Zappella. Random spanning trees and the prediction of weighted graphs. *Journal of Machine Learning Research*, 14, pp. 1251–1284, 2013.
- [8] N. Cesa-Bianchi, C. Gentile, F. Vitale, and G. Zappella. Active learning on trees and graphs. In *Proceedings of the 23rd Conference on Learning Theory (23rd COLT)*, pages 320–332, 2010.
- [9] A. Demiriz, K. Bennett, and M.J. Embrechts. Semi-supervised clustering using genetic algorithms. In *In Artificial Neural Networks in Engineering (ANNIE-99)*, pages 809–814, 1999.

- [10] C. Gentile, M. Herbster, and S. Pasteris. Online similarity prediction of networked data from known and unknown graphs. In *Proceedings of the 23rd Conference on Learning Theory (26th COLT)*, 2013.
- [11] D. Helmbold and M. Warmuth. On weak learning, *J. Comput. System Sci.*, 50, pp. 551–573, 1995.
- [12] M. Herbster. Exploiting cluster-structure to predict the labeling of a graph. In *Proceedings of the 19th International Conference on Algorithmic Learning Theory*, pages 54–69, 2008.
- [13] M. Herbster and G. Lever. Predicting the labelling of a graph via minimum p-seminorm interpolation. In *Proceedings of the 22nd Annual Conference on Learning Theory (COLT'09)*, 2009.
- [14] M. Herbster and M. Pontil. Prediction on a graph with a perceptron. In *Advances in Neural Information Processing Systems 19*, pages 577–584. MIT Press, Cambridge, MA, 2007.
- [15] M. Herbster, M. Pontil, and S. R. Galeano. Fast prediction on a tree. In *Proc. of the 22nd Annual Conference on Neural Information Processing Systems*, pages 657–664. MIT Press, 2008.
- [16] M. Herbster, G. Lever, and M. Pontil. Online prediction on large diameter graphs. In *Advances in Neural Information Processing Systems (NIPS 22)*, pages 649–656. MIT Press, 2009.
- [17] Karger, D. R. Random sampling in cut, flow, and network design problems. In Proc. STOC, 1994.
- [18] O. Kariv, S. Hakimi. An algorithmic approach to network location problems. ii: The p-medians. *SIAM Journal on Applied Mathematics*, 7:3, pages. 539–560, 1979.
- [19] U. Luxborg. A Tutorial on Spectral Clustering. *Statistics and Computing*, 17:4, pages 1251–1284, 2013.
- [20] R. Lyons and Y. Peres. *Probability on Trees and Networks*. Cambridge University Press, 2012. In preparation. Current version available at <http://mypage.iu.edu/~rdlyons/>.
- [21] A. Maurer. Learning similarity with operator-valued large-margin classifiers. *Journal of Machine Learning Research*, 9:1049–1082, 2008.
- [22] Meila, M. (2007). Comparing clusterings—an information based distance. *Journal of Multivariate Analysis*, 98, pp. 873–895.
- [23] Meila, M. (2012). Local equivalences of distances between clusterings—a geometric perspective. *Machine Learning*, 86/3, pp. 369–389.
- [24] Mirkin, B. G. (1996). *Mathematical classification and clustering*. Dordrecht: Kluwer Academic.
- [25] S. S. Rangapuram and M. Hein. Constrained 1-spectral clustering. In *Proc. 15th International Conference on Artificial Intelligence and Statistics, AISTATS*, 2012.
- [26] Rajaraman, A.; Ullman, J. (2010). *Mining of Massive Datasets*.
- [27] Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, 66, 846–850.
- [28] F. Vitale, N. Cesa-Bianchi, C. Gentile, and G. Zappella. See the tree through the lines: The shazoo algorithm. In *NIPS*, pages 1584–1592, 2011.
- [29] M. K. Warmuth. Winnowing subspaces. In *Proceedings of the 24th International Conference on Machine Learning*, pages 999–1006. ACM, 2007.
- [30] E. P. Xing, A. Y. Ng, M. I. Jordan, and J. Russell S. Distance metric learning with application to clustering with side-information. In *NIPS*, pages 505–512, 2002.

A Proofs

A.1 Proof of Theorem 1

The following two lemmas are an immediate consequence of the triangle inequality for DIST.

Lemma 5. *Let $a, b \in [0, 1]$ be such that $a + b \geq 3/2$, and sets U, W, X, Y, Z satisfy*

1. $\text{DIST}(U, W) \leq 1 - a$;
2. $\text{DIST}(W, X) \leq 1 - b$;
3. $\text{DIST}(U, Y) \leq 1 - a$;
4. $\text{DIST}(Z, X) = 1$.

Then $\text{DIST}(Y, Z) \geq 1 - b$.

Proof. We can write

$$\begin{aligned} 1 &= \text{DIST}(Z, X) \\ &\leq \text{DIST}(Z, Y) + \text{DIST}(Y, U) + \text{DIST}(U, W) + \text{DIST}(W, X) \\ &\leq \text{DIST}(Z, Y) + 1 - a + 1 - a + 1 - b, \end{aligned}$$

so that

$$\text{DIST}(Z, Y) \geq 2a + b - 2 \geq 1 - b,$$

the last inequality using the assumption $a + b \geq 3/2$. This concludes the proof. \square

Lemma 6. *Let $a, b \in [0, 1]$ be such that $2b \geq 1 + a$, and sets X, Y, Z satisfy*

1. $\text{DIST}(X, Y) \leq 1 - b$;
2. $\text{DIST}(Y, Z) \geq 1 - a$.

Then $\text{DIST}(X, Z) \geq 1 - b$.

Proof. We can write

$$1 - a \leq \text{DIST}(Y, Z) \leq \text{DIST}(Y, X) + \text{DIST}(X, Z) \leq 1 - b + \text{DIST}(X, Z),$$

so that

$$\text{DIST}(X, Z) \geq b - a \geq 1 - b,$$

the last inequality deriving from $2b \geq 1 + a$. This concludes the proof. \square

With these two simple lemmas handy, we are now ready to analyze RGCA. The reader is compelled to refer to Algorithm 1 for notation. In what follows, $a \in [0, 1]$ is RGCA's distance parameter, and $b \in [0, 1]$ is a constant such that the two conditions on a and b required by Lemmas 5 and 6 simultaneously hold. It is easy to see that these conditions are equivalent to⁴

$$a \geq \frac{2}{3}, \quad b \geq \frac{1+a}{2}. \quad (3)$$

The following definition will be useful.

Definition 7. *A b -anomaly in the similarity graph (V, \mathcal{P}) is a vertex $v \in V$ for which $\text{DIST}(D_{\mu_{\mathcal{D}}(v)}, \Gamma(v)) \geq 1 - b$, for some constant $b \in [0, 1]$ satisfying (3). We denote by Λ_b the set of all anomalies. A centered round of RGCA is any $t \leq \ell$ in which $N_t(\alpha_t) \not\subseteq \Lambda_b$. We denote by Ω_b the set of all centered rounds. A centered label is any class $i \in \{1, \dots, k\}$ such that $D_i \not\subseteq \Lambda_b$. We denote by Δ_b the set of all centered labels.*

Lemma 8. *For any round $t \leq \ell$, there exists a class $j \in \{1, \dots, k\}$ such that for every vertex $v \in N_t(\alpha_t) \setminus \Lambda_b$ we have $\mu_{\mathcal{D}}(v) = j$.*

⁴ For instance, we may set $a = 2/3$ and $b = 5/6$.

Proof. Suppose, for the sake of contradiction, that we have $v, w \in N_t(\alpha_t)$ with $v, w \notin \Lambda_b$ and $\mu_{\mathcal{D}}(v) \neq \mu_{\mathcal{D}}(w)$. Define $U := \Gamma(\alpha_t)$, $W := \Gamma(v)$, $X := D_{\mu_{\mathcal{D}}(v)}$, $Y := \Gamma(w)$, and $Z := D_{\mu_{\mathcal{D}}(w)}$. Since $v, w \in N_t(\alpha_t)$, by the way graph (V, \mathcal{Q}) is constructed, we have both $\text{DIST}(U, W) \leq 1 - a$ and $\text{DIST}(U, Y) \leq 1 - a$. Moreover, since $v \notin \Lambda_b$ we have $\text{DIST}(W, X) < 1 - b$. Also, $\mu_{\mathcal{D}}(v) \neq \mu_{\mathcal{D}}(w)$ implies $\text{DIST}(Z, X) = 1$. We are therefore in a position to apply Lemma 5 verbatim, from which we have $\text{DIST}(Y, Z) \geq 1 - b$, i.e., $w \in \Lambda_b$. This is a contradiction, which implies the claimed result. \square

Lemma 8 allows us to make the following definition.

Definition 9. Given a centered round $t \in \Omega_b$, we define $\gamma(t)$ to be the unique class j such that for every vertex $v \in N_t(\alpha_t) \setminus \Lambda_b$ we have $\mu_{\mathcal{D}}(v) = j$.

Lemma 10. For any round $t \leq \ell$ and vertices $v, w \in A_t$ with $v \notin \Lambda_b$, $w \notin N_t(v)$ and $\mu_{\mathcal{D}}(v) = \mu_{\mathcal{D}}(w)$ we have $w \in \Lambda_b$.

Proof. Define $X := D_{\mu_{\mathcal{D}}(v)}$, $Y := \Gamma(v)$ and $Z := \Gamma(w)$. Since $v \notin \Lambda_b$ we have $\text{DIST}(X, Y) \leq 1 - b$. Moreover, $w \notin N_t(v)$ implies $\text{DIST}(Y, Z) \geq 1 - a$. By Lemma 6 we immediately have $\text{DIST}(X, Z) \geq 1 - b$. But since $X = D_{\mu_{\mathcal{D}}(v)} = D_{\mu_{\mathcal{D}}(w)}$, this equivalently establishes that $w \in \Lambda_b$. \square

Lemma 11. For any centered round $t \in \Omega_b$, any vertex $v \in N_t(\alpha_t) \setminus \Lambda_b$, and any vertex $w \in N_t(\alpha_t) \setminus N_t(v)$, we have $w \in \Lambda_b$.

Proof. If $\mu_{\mathcal{D}}(w) \neq \mu_{\mathcal{D}}(v)$ then by Lemma 8 we must have $w \in \Lambda_b$, so we are done. On the other hand, if $\mu_{\mathcal{D}}(w) = \mu_{\mathcal{D}}(v)$, we have $v \notin \Lambda_b$, $w \notin N_t(v)$ and $\mu_{\mathcal{D}}(v) = \mu_{\mathcal{D}}(w)$ which implies, by Lemma 10, that $w \in \Lambda_b$. \square

Lemma 12. For any centered round $t \in \Omega_b$, we have $|(A_{t+1} \cap D_{\gamma(t)}) \setminus \Lambda_b| \leq |C_t \cap \Lambda_b|$.

Proof. Since $t \in \Omega_b$ there must exist a vertex $v \in N_t(\alpha_t)$ with $v \notin \Lambda_b$, so let us consider such a v . Note that by the way the algorithm works, we have $|N_t(\alpha_t)| \geq |N_t(v)|$, so that $|N_t(v) \setminus N_t(\alpha_t)| \leq |N_t(\alpha_t) \setminus N_t(v)|$. Next, by Lemma 11 we have $N_t(\alpha_t) \setminus N_t(v) \subseteq \Lambda_b$, hence $N_t(\alpha_t) \setminus N_t(v) \subseteq N_t(\alpha_t) \cap \Lambda_b$ and, consequently, $|N_t(\alpha_t) \setminus N_t(v)| \leq |N_t(\alpha_t) \cap \Lambda_b|$. Recalling that $C_t = N_t(\alpha_t)$, we have therefore obtained

$$|N_t(v) \setminus N_t(\alpha_t)| \leq |N_t(\alpha_t) \setminus N_t(v)| \leq |N_t(\alpha_t) \cap \Lambda_b| = |C_t \cap \Lambda_b|. \quad (4)$$

Now suppose we have some vertex $w \in (A_{t+1} \cap D_{\gamma(t)}) \setminus \Lambda_b$. For the sake of contradiction, let us assume that $w \notin N_t(v)$. Then $w \in A_{t+1}$ implies $w \in A_t$ which, combined with Lemma 10 together with the fact that $\mu_{\mathcal{D}}(v) = \gamma(t) = \mu_{\mathcal{D}}(w)$, implies that $w \in \Lambda_b$, which is a contradiction. Hence we must have $w \in N_t(v)$. Moreover, since $w \in A_{t+1}$ we must have $w \notin N_t(\alpha_t)$. We have hence shown that $w \in N_t(v) \setminus N_t(\alpha_t)$, implying that $|(A_{t+1} \cap D_{\gamma(t)}) \setminus \Lambda_b| \leq |N_t(v) \setminus N_t(\alpha_t)|$. Combining with (4) concludes the proof. \square

We now turn to considering centered labels.

Lemma 13. For any centered label $i \in \Delta_b$ there exists some round $t \leq \ell$ such that $\gamma(t) = i$.

Proof. Since i is a centred label, pick $v \in D_i \setminus \Lambda_b$. Further, since C_1, C_2, \dots, C_ℓ is a partition of \mathcal{V} , choose t such that $v \in C_t$. Now, since $v \in C_t \setminus \Lambda_b$ we have that $t \in \Omega_b$ and, by Lemma 8, that $\gamma(t) = \mu_{\mathcal{D}}(v) = i$. \square

Lemma 13 allows us to make the following definition.

Definition 14. Given a centered label $i \in \Delta_b$, we define $\psi(i) := \min\{t : \gamma(t) = i\}$.

Lemma 15. For any centered label $i \in \Delta_b$, we have $D_i \setminus \Lambda_b \subseteq A_{\psi(i)}$.

Proof. Suppose, for contradiction, that there exists some $v \in D_i \setminus \Lambda_b$ with $v \notin A_{\psi(i)}$. Then, by definition of $A_{\psi(i)}$ there exists some round $t^\circ < \psi(i)$ with $v \in C_{t^\circ}$. As $v \notin \Lambda_b$ we have $t^\circ \in \Omega_b$ and, by Lemma 8, that $\mu_{\mathcal{D}}(v) = \gamma(t^\circ)$. Hence $\gamma(t^\circ) = \mu_{\mathcal{D}}(v) = i$ which, due to the condition $t^\circ < \psi(i)$, contradicts the fact that $\psi(i) := \min\{t : \gamma(t) = i\}$. \square

Lemma 16. For any centred label $i \in \Delta_b$ we have $|D_i \setminus C_{\psi(i)}| \leq |D_i \cap \Lambda_b| + |C_{\psi(i)} \cap \Lambda_b|$.

Proof. Suppose we have some $v \in D_i \setminus C_{\psi(i)}$, and let us separate the two cases: (i) $v \notin \Lambda_b$ and, (ii) $v \in \Lambda_b$.

Case (i). Since $v \in D_i \setminus \Lambda_b$ we have, by Lemma 15, that $v \in A_{\psi(i)}$. Since $v \notin C_{\psi(i)}$ this implies that $v \in A_{\psi(i)+1}$. Notice that $\gamma(\psi(i)) = i$ so $D_i = D_{\gamma(\psi(i))}$ and hence $v \in (A_{\psi(i)+1} \cap D_{\gamma(\psi(i))}) \setminus \Lambda_b$. By Lemma 12 the number of such vertices v is hence upper bounded by $|C_{\psi(i)} \cap \Lambda_b|$.

Case (ii). In this case, we simply have that $v \in D_i \cap \Lambda_b$, so the number of such vertices v is upper bounded by $|D_i \cap \Lambda_b|$.

Putting the two cases together gives us $|D_i \setminus C_{\psi(i)}| \leq |D_i \cap \Lambda_b| + |C_{\psi(i)} \cap \Lambda_b|$, as required. \square

Having established the main building blocks of the behavior of RGCA, we now turn to quantifying the resulting connection between ER and HA. To this effect, we start off by defining a natural map Υ associated with the clustering $\{C_1, \dots, C_\ell\}$ generated by RGCA, along with a corresponding accuracy measure.

Definition 17. *The map $\Upsilon : \{D_1, \dots, D_k\} \rightarrow \{C_1, \dots, C_\ell\}$ is defined as follows:*

$$\Upsilon(D_i) = \begin{cases} C_{\psi(i)} & \text{if } i \in \Delta_b \\ \emptyset & \text{if } i \notin \Delta_b \end{cases}$$

Moreover, let $\mathcal{M}(\Upsilon) := \sum_{i=1}^k |D_i \setminus \Upsilon(D_i)|$.

We have the following lemma.

Lemma 18. $\mathcal{M}(\Upsilon) \leq 2|\Lambda_b|$.

Proof. For $i \notin \Delta_b$ we have $D_i \subseteq \Lambda_b$ and $\Upsilon(D_i) = \emptyset$ so that

$$|D_i \setminus \Upsilon(D_i)| = |D_i| = |D_i \cap \Lambda_b| = |D_i \cap \Lambda_b| + |\emptyset| = |D_i \cap \Lambda_b| + |\Upsilon(D_i) \cap \Lambda_b|.$$

On the other hand, for $i \in \Delta_b$ we have $\Upsilon(D_i) = C_{\psi(i)}$ so that, by Lemma 16, we can write

$$|D_i \setminus \Upsilon(D_i)| \leq |D_i \cap \Lambda_b| + |\Upsilon(D_i) \cap \Lambda_b|.$$

Hence, in both cases, for all $i \in \{1, \dots, k\}$ we have

$$|D_i \setminus \Upsilon(D_i)| \leq |D_i \cap \Lambda_b| + |\Upsilon(D_i) \cap \Lambda_b|,$$

implying

$$\mathcal{M}(\Upsilon) \leq \sum_{i=1}^k (|D_i \cap \Lambda_b| + |\Upsilon(D_i) \cap \Lambda_b|). \quad (5)$$

Now, both $\{D_1, \dots, D_k\}$ and $\{\Upsilon(D_1), \dots, \Upsilon(D_k)\}$ are a partition of V , implying

$$|\Lambda_b| = \sum_{i=1}^k |D_i \cap \Lambda_b| = \sum_{i=1}^k |\Upsilon(D_i) \cap \Lambda_b|.$$

Plugging back into (5) yields the claimed result. \square

Next, observe that, by its very definition, $\text{HA}(\mathcal{P}, \mathcal{D})$ can be rewritten as

$$\text{HA}(\mathcal{P}, \mathcal{D}) = \sum_{v \in V} |(D_{\mu_{\mathcal{D}}(v)} \setminus \Gamma(v)) \cup (\Gamma(v) \setminus D_{\mu_{\mathcal{D}}(v)})|. \quad (6)$$

Lemma 19. *We have $\text{HA}(\mathcal{P}, \mathcal{D}) \geq (1 - b) \sum_{i=1}^k d_i |D_i \cap \Lambda_b|$.*

Proof. Fix class $i \in \{1, \dots, k\}$ and vertex $v \in D_i \cap \Lambda_b$. Then $v \in \Lambda_b$ implies $\text{DIST}(D_{\mu_{\mathcal{D}}(v)}, \Gamma(v)) \geq 1 - b$, which in turn yields

$$|(D_{\mu_{\mathcal{D}}(v)} \setminus \Gamma(v)) \cup (\Gamma(v) \setminus D_{\mu_{\mathcal{D}}(v)})| \geq (1 - b)d_i,$$

thereby concluding that for all fixed i

$$\sum_{v \in D_i \cap \Lambda_b} |(D_{\mu_{\mathcal{D}}(v)} \setminus \Gamma(v)) \cup (\Gamma(v) \setminus D_{\mu_{\mathcal{D}}(v)})| \geq (1 - b) |D_i \cap \Lambda_b| d_i.$$

Since $\Lambda_b = \bigcup_{i=1}^k (D_i \cap \Lambda_b)$, being the sets $D_i \cap \Lambda_b$, $i = 1, \dots, k$, pairwise disjoint, we can write

$$\begin{aligned} \sum_{v \in \Lambda_b} |(D_{\mu_{\mathcal{D}}(v)} \setminus \Gamma(v)) \cup (\Gamma(v) \setminus D_{\mu_{\mathcal{D}}(v)})| &= \sum_{i=1}^k \sum_{v \in D_i \cap \Lambda_b} |(D_{\mu_{\mathcal{D}}(v)} \setminus \Gamma(v)) \cup (\Gamma(v) \setminus D_{\mu_{\mathcal{D}}(v)})| \\ &\geq (1-b) \sum_{i=1}^k |D_i \cap \Lambda_b| d_i. \end{aligned}$$

Thus, from (6), and the fact that $\Lambda_b \subseteq \mathcal{V}$ the result immediately follows. \square

Lemma 20. *The number $|\Lambda_b|$ of b -anomalies can be upper bounded as*

$$|\Lambda_b| \leq \min_{j=1, \dots, k} \left(\frac{1}{d_j(1-b)} \text{HA}(\mathcal{P}, \mathcal{D}) + \sum_{i=1}^{j-1} d_i \right).$$

Proof. For any $j = 1, \dots, k$ we can write

$$|\Lambda_b| = \sum_{i=1}^k |D_i \cap \Lambda_b| = \sum_{i=1}^{j-1} |D_i \cap \Lambda_b| + \sum_{i=j}^k |D_i \cap \Lambda_b| \leq \sum_{i=1}^{j-1} d_i + \sum_{i=j}^k |D_i \cap \Lambda_b|$$

so all that is left to prove is that the last sum in the right-hand side is at most $\frac{1}{d_j(1-b)} \text{HA}(\mathcal{P}, \mathcal{D})$.

Since, for all classes i such that $i \geq j$, we have $d_i \geq d_j$, we can write

$$\sum_{i=j}^k |D_i \cap \Lambda_b| \leq \sum_{i=j}^k \frac{d_i}{d_j} |D_i \cap \Lambda_b| \leq \frac{1}{d_j} \sum_{i=1}^k d_i |D_i \cap \Lambda_b| \leq \frac{1}{d_j(1-b)} \text{HA}(\mathcal{P}, \mathcal{D}),$$

where the last inequality derives from Lemma 19. This concludes the proof. \square

We are now ready to combine to above lemmas into the proof of Theorem 1.

Proof. (Theorem 1) Direct from Lemmas 18 and 20 we have

$$\mathcal{M}(\Upsilon) \leq \min_{j=1, \dots, k} \left(\frac{2}{d_j(1-b)} \text{HA}(\mathcal{P}, \mathcal{D}) + \sum_{i=1}^{j-1} d_i \right).$$

We then optimize for b by selecting $b = \frac{1+a}{2}$, and then for a by setting $a = 2/3$, so as to fulfil conditions (3). The result follows by the fact that $\text{ER}(\mathcal{C}, \mathcal{D}) \leq \mathcal{M}(\Upsilon)$, for $\text{ER}(\mathcal{C}, \mathcal{D})$ is a minimum over all possible cluster maps $\mathcal{D} \rightarrow \mathcal{C}$, while Υ is just the one in Definition 17. \square

A.2 Proof of Theorem 2

Proof. For ease of proof, we assume that d_j is even for all j (adapting the proof to the general case is trivial). We consider two cases:

1. $\sigma \geq \frac{1}{2} \sum_{j=1}^k d_j^2$;
2. $\sigma < \frac{1}{2} \sum_{j=1}^k d_j^2$.

For the first case we choose, for every $j = 1, \dots, k$, sets P_j^+ and P_j^- such that $|P_j^+| = |P_j^-| = d_j/2$ and $P_j^+ \cup P_j^- = D_j$. We then construct the similarity graph $(\mathcal{V}, E_{\mathcal{P}})$, where clustering \mathcal{P} is made up of the $2k$ clusters $\{P_j^+ : j = 1, \dots, k\} \cup \{P_j^- : j = 1, \dots, k\}$. Since the algorithm is consistent, we must have $\mathcal{C} = \mathcal{P}$. Now, let f be an injection from \mathcal{D} to \mathcal{C} , and consider any $j = 1, \dots, k$. If $f(D_j) \in \{P_j^+, P_j^-\}$ then we have $|D_j \setminus f(D_j)| = d_j/2$, and otherwise $|D_j \setminus f(D_j)| = d_j$, so that

$$\sum_{j=1}^k |D_j \setminus f(D_j)| \geq \frac{1}{2} \sum_{j=1}^k d_j = n/2.$$

Since f is arbitrary, this shows that $\text{ER}(\mathcal{C}, \mathcal{D}) \geq \frac{\sigma}{2}$. Moreover, we observe that the only incorrect similarity/dissimilarity predictions of \mathcal{P} with respect to \mathcal{D} are those between P_j^+ and P_j^- , for every j , which gives us $2|P_j^+| \cdot |P_j^-| = d_j^2/2$ incorrect predictions for every j . This implies that $\text{HA}(\mathcal{P}, \mathcal{D}) = \sum_{j=1}^k d_j^2/2$, which is no greater than σ , thereby completing the proof for the first case.

We now turn to the second case. Let $j^\circ \in \{1, \dots, k\}$ be such that

$$\frac{1}{2} \sum_{i=1}^{j^\circ-1} d_i^2 \leq \sigma < \frac{1}{2} \sum_{i=1}^{j^\circ} d_i^2,$$

and $\omega := \sigma - \frac{1}{2} \sum_{i=1}^{j^\circ-1} d_i^2$. Notice that $\omega \leq d_{j^\circ}^2/2$. We choose, for every $j < j^\circ$, sets P_j^+ and P_j^- such that $|P_j^+| = |P_j^-| = d_j/2$ and $P_j^+ \cup P_j^- = D_j$. Let $c = \lfloor \omega/2d_{j^\circ} \rfloor$, and note that $c \leq d_{j^\circ}/4 < d_{j^\circ}/2$. We can hence define subsets $X, Y \subseteq D_{j^\circ}$ such that $|X| = c$, $X \cup Y = D_{j^\circ}$ and $X \cap Y = \emptyset$.

We construct the similarity graph $(\mathcal{V}, E_{\mathcal{P}})$, where clustering \mathcal{P} is made up of the $k + j^\circ$ clusters

$$\{P_j^+ : j = 1, \dots, j^\circ - 1\} \cup \{P_j^- : j = 1, \dots, j^\circ - 1\} \cup \{X, Y\} \cup \{D_j : j > j^\circ\}.$$

Again, since the algorithm is consistent, we must have $\mathcal{C} = \mathcal{P}$. As before, let f be an arbitrary injection from \mathcal{D} to \mathcal{C} , and consider any $j < j^\circ$. Then if $f(D_j) \in \{P_j^+, P_j^-\}$ we have $|D_j \setminus f(D_j)| = d_j/2$, otherwise $|D_j \setminus f(D_j)| = d_j$, so that $|D_j \setminus f(D_j)| \geq d_j/2$ holds for any $j < j^\circ$. Further, if $f(D_{j^\circ}) = X$ then $|D_{j^\circ} \setminus f(D_{j^\circ})| = d_{j^\circ} - c$, if $f(D_{j^\circ}) = Y$ then $|D_{j^\circ} \setminus f(D_{j^\circ})| = c$, and otherwise $|D_{j^\circ} \setminus f(D_{j^\circ})| = d_{j^\circ}$. In any case, since $c < d_{j^\circ}/2$, we have $|D_{j^\circ} \setminus f(D_{j^\circ})| \geq c$. This allows us to conclude that

$$\begin{aligned} \text{ER}(\mathcal{C}, \mathcal{D}) &= \sum_{j=1}^k |D_j \setminus f(D_j)| \\ &\geq c + \frac{1}{2} \sum_{j=1}^{j^\circ-1} d_j \\ &= \lfloor \omega/2d_{j^\circ} \rfloor + \frac{1}{2} \sum_{j=1}^{j^\circ-1} d_j \\ &\geq \frac{\omega}{2d_{j^\circ}} - 1 + \frac{1}{2} \sum_{j=1}^{j^\circ-1} d_j \\ &= \frac{\sigma}{2d_{j^\circ}} - 1 - \frac{1}{4d_{j^\circ}^2} \sum_{j=1}^{j^\circ-1} d_j^2 + \frac{1}{2} \sum_{j=1}^{j^\circ-1} d_j \\ &= \frac{\sigma}{2d_{j^\circ}} - 1 + \frac{1}{2} \sum_{j=1}^{j^\circ-1} d_j \left(1 - \frac{d_j}{2d_{j^\circ}}\right) \\ &\geq \frac{\sigma}{2d_{j^\circ}} - 1 + \frac{1}{4} \sum_{j=1}^{j^\circ-1} d_j. \end{aligned}$$

Finally, notice that the only incorrect similarity/dissimilarity predictions of \mathcal{P} with respect to \mathcal{D} are those between P_j^+ and P_j^- , for every $j < j^\circ$, and those between X and Y , which gives us $2|P_j^+| \cdot |P_j^-| = d_j^2/2$ incorrect predictions for every $j < j^\circ$, and an additional $2|X| \cdot |Y| = 2c(d_{j^\circ} - c) \leq 2cd_{j^\circ} \leq \omega$ incorrect predictions between X and Y . This implies that

$$\text{HA}(\mathcal{P}, \mathcal{D}) \leq \omega + \sum_{j=1}^{j^\circ-1} d_j^2/8$$

which is in turn bounded from above by σ . This completes the proof for the second case. \square

The following simple lemma is of preliminary importance for the proof of Theorem 3.

Lemma 21. *Let $H = (V, E)$ be an Erdos-Renyi $G(n, p)$ graph. For each subgraph $H'(V', E') \subseteq H$ with $n' = |V'|$ nodes, when $p = \frac{\lambda \log n'}{n'}$ the following separation property holds: As n' approaches infinity, the expected number z of isolated vertices in G' equals $(n')^{1-\lambda}$. Furthermore, in the special case when $n' = \frac{1}{p}$, we always have $z \geq \frac{1}{pe}$.*

Proof. In order to prove these properties, it suffices to observe that, given any node in V' , the probability that it is isolated in G' is equal to $(1-p)^{n'-1}$, which in turn is equal to $e^{-\lambda \log n'} = (n')^{-\lambda}$ as n' approaches infinity. Hence we have $z = (n')^{1-\lambda}$. By a similar argument, it is immediate to verify that in the case when $n' = \frac{1}{p}$ we have $z = \frac{1}{p}(1-p)^{\frac{1}{p}-1}$ which is never smaller than $\frac{1}{ep}$. \square

A.3 Proof of Theorem 3

Proof. Let $G' = (V, E')$ denote the undirected graph whose edge set E' is made up of all pairs of vertices drawn in S . Since S is drawn uniformly at random, G' turns out to be an Erdos-Renyi graph $G'(n, p)$.

Setting $\lambda = 2$ in Lemma 21, we have that for all clusters $C \in \mathcal{C}$ such that $\frac{2 \log |C|}{|C|} \leq p$, cluster C can be *completely* detected by SACA (line 3 in Algorithm 2), with probability at least $\frac{1}{|C|}$. Hence, the expected number of misclassification errors made when detecting such clusters is upper bounded by 1 per cluster. In order to satisfy the assumption $\frac{2 \log |C|}{|C|} \leq p$, the size of these clusters must be equal to a value $\tau = \Omega(\rho \log \rho)$, where we set $\rho = \frac{1}{p}$.

Finally, we can conclude the proof observing that the total number of misclassification errors is bounded in expectation by the sum of the following two quantities: (i) the number of clusters larger than τ , which in turn is bounded by k , and (ii) the total number of nodes belonging to the clusters smaller or equal to τ , which in turn is bounded by $k\tau$:

$$\mathbb{E}[\text{ER}(\mathcal{C}_{\text{final}}, Y)] = \mathcal{O}(k(1 + \tau)) = \mathcal{O}(k\rho \log \rho), \quad (7)$$

thereby concluding the proof. \square

A.4 Proof of Theorem 4

Proof. As in the proof of Theorem 3, we denote by $G' = (V, E')$ the undirected graph whose edge set E' is made up of all pairs of vertices drawn in the training set S . Since S is drawn uniformly at random, G' turns out to be an Erdos-Renyi graph.

The basic idea in this proof is to construct a collection \mathcal{H} of z disjoint subsets of V , call them H_1, H_2, \dots, H_z , and, for all $j \in \{1, \dots, z\}$, to *randomly* label all nodes of each subset H_j using only a pair of classes of $\{1, \dots, k\}$. These z pairs of classes must be distinct and disjoint. The random labeling is accomplished in such a way that no algorithm can exploit the training set nor the information carried by G to guess how each H_j is labeled, while we always guarantee $\Phi_G^R(Y) \leq b$. More specifically, \mathcal{H} is created so as to satisfy the following three properties:

Property (i) For all $j = 1, \dots, z$, no pair of nodes in H_j are connected by an edge in the training graph representation G' , i.e. for each pair of nodes $u, v \in H_j$, we have $(u, v) \notin S$.

Property (ii) Consider *any* possible vertex labeling from $\{1, \dots, k\}$ of all sets in \mathcal{H} that uses at most $k-1$ classes, say $1, \dots, k-1$. Then if we assign label k (which is never used for vertices in the sets of H_j , to all the remaining nodes in V , i.e. those in $V \setminus \cup_{j=1}^z H_j$, we can *always* ensure that $\Phi_B^R \leq b$.

Property (iii) For all $j = 1, \dots, z$, we have that the expected size of H_j (over the random draw of the training set S) is larger than $\frac{n^2}{2me}$, if $m > \frac{n}{2}$, where e is the base of natural logarithms, while it is $\Theta(b)$ if $m \leq \frac{n}{2}$.

Figure 1 provides a pictorial explanation of the randomized labeling strategy we are going to describe.

We now describe in detail the randomized labeling strategy (a randomized similarity matrix Y representing a clustering with k clusters), and derive a lower bound for $\mathbb{E}_Y[\text{ER}(\mathcal{C}, Y)]$ when \mathcal{H} satisfies all of the above properties.

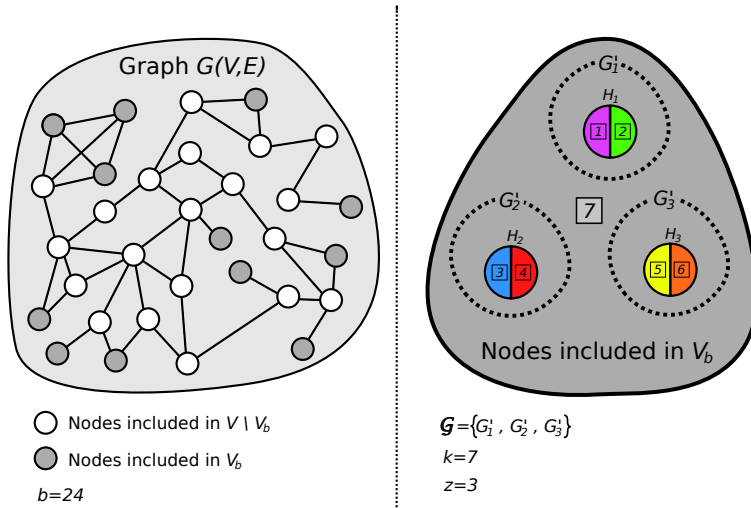


Figure 1: Illustration of the randomized labeling that achieves the lower bound in Theorem 4. **Left:** The side information graph $G = (V, E)$. Here $b = 24$. V_b is therefore made up of all $\frac{24}{2} = 12$ grey nodes in the picture, which are the 12 nodes with smallest $r(v)$ values among all nodes $v \in V$. Since $R(12) < 24$, if all white nodes are assigned to the same class we can ensure $\Phi_G^R < b = 2 \cdot 12 = 24$, independent of the chosen labels for the grey nodes. **Right:** The grey area includes all the nodes of V_b . For the depicted graph, we have $z = \lfloor \frac{2m}{n^2} |V_b| \rfloor = 3$ and $k = 7$. In this case we thus have $\lfloor \frac{k-1}{2} \rfloor = 3$. \mathcal{G} is the collection of the 3 vertex-disjoint subgraphs G'_1, G'_2 and G'_3 . The node set size of each of these subgraphs is equal to $\lfloor \frac{n^2}{2m} \rfloor$. The subsets of isolated vertices in these 3 subgraphs are H_1, H_2 , and H_3 , which are depicted in this figure by the bicoloured circles. Each color represents a class. For each j , the expected size of H_j must be linear in the size of the node set of G'_j . For $j = 1, 2, 3$, set H_j is labeled by selecting uniformly at random a class between the two classes (or colors) $2j - 1$ and $2j$. All the remaining nodes in the grey area of this picture (together with the white nodes in the picture on the left) are given the same class 7. Hence, for each pair of nodes u and v both belonging to H_j for some j , we must have $(u, v) \notin S$. On the contrary, for each pair u and v with $u \in H_j$, for some j , and $v \notin H_j$, we must have $y_{u,v} = 0$. Neither the information of the training set nor the graph topology of G can be used to predict how the nodes in H_1, H_2 , and H_3 , are labeled. In fact, *any* algorithm will make $\frac{1}{2}$ mistakes in expectation over the randomized labeling on these nodes. On the other hand, it holds by construction that $\Phi_G^R \leq b$.

Let $z \leq \lfloor \frac{k-1}{2} \rfloor$. Once we constructed such a collection \mathcal{H} of clusters, we associate a distinct pair of classes in $\{1, \dots, k\}$ with each H_j in such a way that all these class pairs are distinct and disjoint. This allows us to always leave one class out (say, class k) for labeling all remaining vertices in V . In particular, we associate with H_j with the class pair $(2j - 1, 2j)$, and then adopt the following randomized strategy:

For all $j = 1, \dots, z$, set H_j is split uniformly at random into two subsets H'_j and H''_j , and we label H'_j by class $2j - 1$ and H''_j by class $2j$. All remaining nodes in $V \setminus \cup_{j=1}^z H_j$ are labeled with class k .

This randomized labeling strategy ensures that, in order to guess the true clustering Y , no learning algorithm can exploit the information provided by S , since for all node pairs (v, w) with $v \in H_j$, for some $j \in \{1, \dots, z\}$, one of two cases hold:

Case (a): $w \in H_j$, which implies that $(v, w) \notin S$, because of Property (i). We have therefore no training set information related to the similarity of nodes laying in the same set H_j .

Case (b): $w \notin H_j$. In this case, whenever $(v, w) \in S$, we *always* have $y_{v,w} = 0$, and this information cannot be exploited to guess the randomized labeling of H_j .

In short, no training information can be exploited to guess how each set H_j is split into the two subsets H'_j and H''_j . Furthermore, the side information graph G cannot be exploited because the randomized labeling is completely independent of the G 's topology, while Property (ii) ensures that \mathcal{H} is selected in such a way that we always have $\Phi_G^R \leq b$. This entails that any clustering algorithm

will incur an expected number of misclassification errors proportional to

$$\sum_{j=1}^z |H_j| = \Omega \left(\min \left\{ \frac{n^2}{m} z, n \right\} \right),$$

the latter equality deriving from Property (iii).

We now turn to describing the detailed construction of \mathcal{H} . We first need to find $V_b \subset V$ such that $|V_b| = \lfloor \frac{b}{2} \rfloor$, and for any labeling of V such that all nodes contained in $V \setminus V_b$ are uniformly assigned to the same class, we always have $\Phi_G^R \leq b$. All sets H_1, \dots, H_z are subsets of V_b , so as to satisfy Property (ii). Note that the assumption $b \geq 4$ is used here to ensure $|V_b| > 1$. Thereafter, we will explain how to select the z subsets satisfying Property (i), and show that their size is bounded from below as required by Property (iii). This will lead to the claimed lower bound.

Satisfaction of property (ii). Let $r_{v,w}$ be the effective resistance between nodes v and w in $G = (V, E)$, and $r(v) = \sum_{w: (v,w) \in E} r_{v,w}$. Moreover, given any integer $h \leq n$, let $R(h) = \min_{\{v_1, v_2, \dots, v_h\} \subseteq V} \sum_{\ell=1}^h r(v_\ell)$. We have

$$R(h) \leq \frac{h}{n} \sum_{v \in V} r(v) = \frac{2h}{n} \sum_{(v,w) \in E} r_{v,w} = \frac{2h(n-1)}{n} < 2h.$$

Let now V_b be the subset of V containing the $\lfloor \frac{b}{2} \rfloor$ nodes $v_1, v_2, \dots, v_{\lfloor b/2 \rfloor}$ achieving the smallest values of $r(v)$ over all $v \in V$. Hence we must have $R(\lfloor \frac{b}{2} \rfloor) \leq b$, which implies that for any vertex labeling such that all nodes in $V \setminus V_b$ are labeled with the same class, we have $\Phi_G^R \leq b$. As anticipated, we will construct \mathcal{H} using only subsets from V_b , and we assign to all nodes in $V \setminus V_b$ the same class, thereby fulfilling Property (ii).

Definition of z . Let

$$z = \min \left\{ f(b, n, m), \left\lfloor \frac{k-1}{2} \right\rfloor \right\}, \quad \text{where} \quad f(b, n, m) = \max \left\{ \left\lfloor \frac{bm}{n^2} \right\rfloor, 1 \right\}.$$

Satisfaction of Property (i).

Let \mathcal{H}' be a collection of disjoint subsets of V_b created as follows. If $\lfloor \frac{b}{2} \rfloor < \lfloor \frac{n^2}{2m} \rfloor$, then $f(b, n, m)$ and z are both equal to 1 and \mathcal{H}' contains only V_b . In all other cases, \mathcal{H}' is generated by selecting uniformly at random z disjoint subsets of V_b such that each node subset contains $\lfloor \frac{n^2}{2m} \rfloor$ nodes (observe that even in the latter case we may have $z = 1$ when $k = 3$ or $1 \leq \lfloor \frac{bm}{n^2} \rfloor < 2$). The collection of subsets $\mathcal{H} = \{H_1, \dots, H_z\}$ is constructed as described next. Let $\mathcal{G} \equiv \{G'_1, G'_2, \dots, G'_z\}$, where G'_j is the subgraph of G' induced by the nodes in the j -th set of \mathcal{H}' . We create z -many disjoint subsets H_1, H_2, \dots, H_z by selecting all vertices that are *isolated* in each graph of \mathcal{G} , and set $\mathcal{H} \equiv \{H_1, H_2, \dots, H_z\}$. Property (i) is therefore satisfied.

Satisfaction of property (iii). By definition of \mathcal{H}' , each graph of the collection \mathcal{G} has $\lfloor \frac{n^2}{2m} \rfloor$ nodes. Using the second part of Lemma 21, we conclude that the expected size of each set in \mathcal{H} is not smaller than $\frac{\lfloor n^2/2m \rfloor}{e}$.

Now, if $\lfloor \frac{b}{2} \rfloor \geq \lfloor \frac{n^2}{2m} \rfloor$, because the set of vertices of each graph in \mathcal{G} is made up of $\lfloor \frac{n^2}{2m} \rfloor$ nodes, we can conclude that the expected number of isolated nodes contained in each graph of the collection \mathcal{G} is linear in its vertex set size. If instead we have $\lfloor \frac{b}{2} \rfloor < \lfloor \frac{n^2}{2m} \rfloor$, by the way $f(b, n, m)$ is defined and because of property (iii), the number of isolated nodes contained in V_b is linear in the size of V_b itself. Furthermore, the size of the (unique) set of nodes contained in \mathcal{H} is linear in b , as well as the expected number of mistakes that can be forced. In fact, the claimed lower bound is $\Omega(b)$ when $\lfloor \frac{b}{2} \rfloor < \lfloor \frac{n^2}{2m} \rfloor$.

Hence the collection of sets \mathcal{S} so generated fulfils at the same time Properties (i), (ii) and (iii).

In order to conclude the proof, we compute our lower bound based on the definition of z and Property (iii). As anticipated, because of the randomized labeling strategy, the expected number of

misclassification errors made by any algorithm is proportional to $\sum_{j=1}^z |H_j| = \Omega \left(\min \left\{ \frac{n^2}{m} z, n \right\} \right)$.
 Plugging in the values of z yields

$$\sum_{j=1}^z |H_j| = \Omega \left(\min \left\{ \frac{\min \left\{ \max \left\{ \lfloor \frac{bm}{n^2} \rfloor, 1 \right\}, \lfloor \frac{k-1}{2} \rfloor \right\}}{m/n^2}, n \right\} \right) = \Omega \left(\min \left\{ \frac{n^2}{m} k, b \right\} \right),$$

and the proof is concluded. □