

An Incremental Reseeding Strategy for Clustering

Xavier Bresson and Huiyi Hu and Thomas Laurent and Arthur Szlam and James von Brecht

Abstract We propose an easy-to-implement and highly parallelizable algorithm for multiway graph partitioning. The algorithm proceeds by alternating three simple routines in an iterative fashion: diffusion, thresholding, and random sampling. We demonstrate experimentally that the proper combination of these ingredients leads to an algorithm that achieves state-of-the-art performance in terms of cluster purity on standard benchmark data sets. We also describe a coarsen, cluster and refine approach similar to [1, 2] that removes an order of magnitude from the runtime of our algorithm while still maintaining competitive accuracy.

1 Introduction

One of the most basic unsupervised learning tasks is to automatically partition data into clusters based on similarity. A standard scenario is that the data are represented as a weighted graph. Data points correspond to vertices on the graph while edges between vertices encode the similarity between data points. Many of the most popular and widely used clustering algorithms, such as spectral clustering, fall into this category. Despite the vast literature on graph-based clustering, the field remains an active area for both theoretical and practical research.

In this work, we propose a resampling-based spectral algorithm for multiway graph partitioning. On graphs that contain reasonably well-balanced clusters of medium scale, the algorithm provides a strong combination of accuracy, efficiency and robustness to noise in the graph construction process. The algorithm is also exceedingly simple, intuitive and trivial to implement. A MATLAB code consists of fewer than 40 lines, for instance (the code is provided in the appendix). It also parallelizes trivially, and can therefore scale gracefully to large numbers of clusters as well as to graphs with large numbers of vertices.

We validate these claims via an extensive experimental evaluation of the algorithm. We also provide a detailed algorithmic comparison using recent clustering algorithms that claim state-of-the-art results. These experiments demonstrate that our algorithm achieves state-of-the-art performance in terms of cluster purity while running faster than the other highly accurate clustering methods (e.g. [3, 4]) that we compare against. We also provide experiments to demon-

Xavier Bresson

School of Computer Science and Engineering, Nanyang Technological University, Singapore, e-mail: xbresson@ntu.edu.sg

Huiyi Hu

Google, e-mail: clarahu@google.com

Thomas Laurent

Department of Mathematics, Loyola Marymount University, e-mail: tlaurent@lmu.edu

Arthur Szlam

Facebook AI Research, New York, e-mail: aszlam@fb.com

James von Brecht

Department of Mathematics, California State University, Long Beach, e-mail: james.vonbrecht@csulb.edu

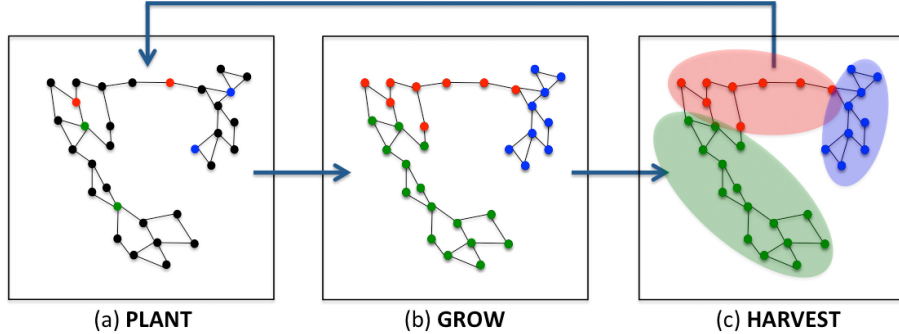


Fig. 1 Illustration of the Incremental Reseeding (INCRES) Algorithm for $R = 3$ clusters. The various colors red, blue, and green identify the clusters. Figure (a): At this stage of the algorithm, $s = 2$ seeds are randomly planted in the clusters computed from the previous iteration. Figure (b): The seeds grow with the random walk operator. Figure (c): A new partition of the graph is obtained and used to plant $s + ds$ seeds into these clusters at the next iteration.

strate the robustness of the algorithm with respect to noise and perturbations in the underlying graph. While many highly accurate algorithms exhibit a sharp decrease in accuracy if the input graph is corrupted by noise, our algorithm remains stable: the accuracy of our algorithm decays slowly and gracefully with increasing levels of noise. These results, when taken together, lead to an algorithm with a quite appealing combination of simplicity, performance and ease in out-of-the-box usage.

2 Description of the Algorithm

The main idea behind our algorithm arises from a well-known and widely used property of the random walk on a graph. Specifically, a random walker started in a low conductance cluster is unlikely to leave that cluster quickly [5]. This fact provides the basis for transductive label propagation methods [6] as well as for “local” clustering methods [7]. In label propagation, for instance, an oracle provides a set of labeled vertices that are propagated along the graph using a random walk matrix or a diffusion matrix. Each unlabeled vertex is then associated to the label which, after being propagated, best represents the given unlabeled vertex.

Our algorithm simply iterates upon this basic idea. Assume that the graph has R well-defined clusters of comparable size and low conductance. If we knew these clusters in advance, we could then select a handful of “seed” vertices in the center of each cluster. We would then expect to obtain good results from a transductive label propagation by using these seeds as labels. In an unsupervised context we cannot, of course, *a-priori* place seeds in the center of each cluster. To overcome this, we instead place a handful of seeds at random. We then apply a random walk matrix or diffusion matrix a few times to propagate these seeds. We finally obtain a temporary clustering by assigning each vertex to the seed which, after propagation, best represents the vertex. We then choose new seeds from these temporary clusters and iterate the process. If the clusters improve then the seeds will likely improve, and vice-versa. This incites a feed-back loop and we get a virtuous cycle. We can then excite the speed and improve the quality of this cycle by gradually drawing more and more seeds throughout the process. We refer to this idea as an *incremental reseeding strategy*. Figure 1 depicts this cyclic process graphically.

2.1 The Basic Algorithm

To formalize these ideas, let $G = (V, W)$ denote a weighted, connected graph on N vertices $V = \{v_1, \dots, v_N\}$ with edge weights $W = \{W_{ij}\}_{i,j=1}^N$ that encode a measure of similarity between each pair (i, j) of vertices. Let D denote the diagonal matrix of (weighted) vertex degrees. The algorithm starts from a random partition $\mathcal{P} = (\mathcal{C}_1, \dots, \mathcal{C}_R)$,

$\mathcal{C}_1 \cup \dots \cup \mathcal{C}_R = V$, $\mathcal{C}_r \cap \mathcal{C}_q = \emptyset$ ($r \neq q$) of the vertices. In other words, each v_i is assigned to one of the R clusters uniformly at random. Let $s = 1$ denote the initial number of seeds. At each of the successive iteration, we update the current partition $\mathcal{P} = (\mathcal{C}_1, \dots, \mathcal{C}_R)$ according to the steps outlined in Algorithm 1. At the beginning of each iteration, the routine PLANT(\mathcal{P}, s) will sample s seeds from each of the R clusters \mathcal{C}_r in the current partition \mathcal{P} uniformly at random. These Rs seeds furnish the temporary labels that the GROW routine then propagates along the graph using a random walk matrix. We initialize GROW with an $N \times R$ matrix U , where $U_{i,r} = 1$ if vertex v_i was drawn from cluster \mathcal{C}_r and $U_{i,r} = 0$ otherwise. We then iteratively apply the random walk transition matrix WD^{-1} to U until each vertex has a nonzero probability of being visited by a random walker, i.e. until each entry $U_{i,r}$ is nonzero. Finally, the routine HARVEST simply assigns each vertex v_i to its most likely cluster, or in other words the cluster for which $U_{i,r}$ is maximal. This produces a new partition \mathcal{P} of the vertices into R clusters. We then increment s to $s + ds$, and use this partition and number of seeds s to initialize PLANT at the beginning of the next iteration. We refer to this overall procedure as the Incremental Reseeding Algorithm (INCRES).

Algorithm 1 INCRES Algorithm

Input: Similarity matrix W , seed increment ds , number of clusters R .
Initialization: $s = 1$, random partition \mathcal{P} .
repeat
 $F \leftarrow \text{PLANT}(\mathcal{P}, s)$
 $U \leftarrow \text{GROW}(F, W)$
 $\mathcal{P} \leftarrow \text{HARVEST}(U)$
 $s \leftarrow s + ds$
until \mathcal{P} converges
Output: \mathcal{P}

function PLANT(\mathcal{P}, s)
 Initialize F as an N -by- R matrix of zeros.
 for $r = 1$ to R **do**
 for $k = 1$ to s **do**
 Draw at random a vertex i in cluster \mathcal{C}_r .
 $F_{i,r} \leftarrow F_{i,r} + 1$
 end for
 end for
 return F .
end function

function GROW(F, W)
 Initialize U as an N -by- R matrix equal to F .
 while $\min_r \min_r U_{i,r} = 0$ **do**
 $U \leftarrow (WD^{-1})U$
 end while
 return U .
end function

function HARVEST(U)
 for $r = 1$ to R **do**
 $\mathcal{C}_r = \{i : U_{i,r} \geq U_{i,q} \text{ for all } q\}$
 end for
 return $\mathcal{P} = (\mathcal{C}_1, \dots, \mathcal{C}_R)$
end function

The overall routine has only a single parameter ds that controls the linear rate at which the number of seeds drawn at each iteration increases. In practice, we select

$$ds = \mathbf{speed} \times 10^{-4} \times \frac{N}{R} \quad (1)$$

for some proportionality constant **speed** between one and ten. By rescaling ds in this way, a constant of proportionality **speed** = 1 corresponds to a total of $s = 0.1N/R$ seeds planted in each cluster after 1000 iterations. Assuming well-balanced clusters of roughly equal size, approximately one-tenth of each cluster is sampled after 1000 iterations. Drawing a significant fraction of each cluster will cause the subsequent clustering to stabilize, leading to eventual convergence of the algorithm. Using around 10% of labels in each cluster is a typical level at which INGRES will stabilize.

The parameter ds therefore represents a “timestep” for INGRES, and the overall algorithm behaves well with respect to this parameter. In general, small increments ds will lead to slower convergence at higher accuracy while larger increments ds will lead to faster convergence but potentially less accurate solutions. Our experiments show that **speed** = 1 works remarkably well for a large variety of data sets. We also provide results for **speed** = 5, which yields faster stabilization with slightly less accuracy, to show that the algorithm is indeed robust and predictable with respect to the choice of this parameter.

The general INGRES framework also proves robust to implementation choices for the three main routines. For instance, in addition to the random walk matrix WD^{-1} , there exists a variety of alternative means to propagate labels along a graph. By-and-large, the overall INGRES strategy does not depend heavily upon the particular implementation of GROW, so long as it realizes the basic idea of label propagation in one form or another. For instance, we have found that replacing the random walk step $U \leftarrow WD^{-1}U$ with a diffusion step $U \leftarrow D^{-1}WU$ or $U \leftarrow D^{-1/2}WD^{-1/2}U$ will give similar results in many circumstances. Occasionally, we have found that utilizing a “personalized Page-Rank” step

$$U \leftarrow \alpha WD^{-1}U + (1 - \alpha)F \quad (2)$$

can give better performance on small data sets that contain a large (relative to the size of the data set) number of clusters. Here the parameter $0 < \alpha < 1$ denotes a length-scale that controls the extent of diffusion and F denotes the input to the GROW routine. A propagation step of the form (2) is also used in Pagerank-NIBBLE [8] and NMFR [3], up to replacing WD^{-1} with $D^{-1/2}WD^{-1/2}$ in the latter case. As another example, choosing to sample with or without replacement in the PLANT routine leads to essentially no significant difference in the resultant clusterings.

2.2 Relation with Other Work

Our methodology relies upon and incorporates number of ideas from transductive learning. In particular, we leverage the notion of label propagation [6]. In the standard label propagation framework, an oracle provides a set of labeled points or vertices. These labeled vertices form either nonzero initial conditions or heat sources for a discrete heat equation on the graph. The second step of the INGRES algorithm (the GROW routine) precisely corresponds with a label propagation of the random labels returned from the first step of the algorithm (the PLANT routine).

The proposed algorithm has a quite intuitive and appealing motivation based on a first principle approach to graph partitioning, in the sense that INGRES combines a simple iterative application of label propagation and thresholding to handle unsupervised learning. At a deeper level, we may also view INGRES as a leading-order approximation of the algorithm from [4] that optimizes the product cut (PCUT) objective. The PCUT algorithm optimizes the “likelihood”

$$\mathcal{L}(\mathcal{P}) := \prod_{r=1}^R \prod_{v_j \in \mathcal{C}_r} \text{prob}_{\mathcal{C}_r}(v_j)$$

of a partition, where the “probability distribution” $\text{prob}_{\mathcal{C}_r}(v_j)$ of a cluster $\mathcal{C}_r \subsetneq V$ results from iterating the personalized Page-Rank step (2) until convergence. The subsequent optimization of $\mathcal{L}(\mathcal{P})$ from [4] then proceeds using a sequence of three routines analogous to the PLANT, GROW and HARVEST routines of INGRES. While the PLANT and HARVEST strategies from [4] barely differ from the INGRES algorithm, the GROW routine

function GROWPCUT(F, W)

$\hat{F} = F \text{diag}(1^T F)^{-1}$
 Solve $(\text{Id} - \alpha W D^{-1}) \hat{U} = (1 - \alpha) \hat{F}$ for \hat{U}

$\tilde{F}_{ik} = F_{ik} / \hat{U}_{ik}$
 Solve $(\text{Id} - \alpha D^{-1} W) \tilde{U} = (1 - \alpha) \tilde{F}$ for \tilde{U}

return $U = \tilde{U} + \log \hat{U}$.

end function

for PCUT requires two label propagation steps performed in series. After normalizing F to make \hat{F} column-stochastic, the solution \hat{U} of the first system comes from iterating

$$U \leftarrow \alpha W D^{-1} U + (1 - \alpha) \hat{F} \quad (3)$$

until convergence. The solution \tilde{U} of the second system similarly comes from an iterative diffusive process

$$U \leftarrow \alpha D^{-1} W U + (1 - \alpha) \tilde{F}. \quad (4)$$

In other words, the GROW function for the PCUT algorithm also performs a type of label propagation on randomly planted seeds, but it does so in quite an expensive way. The dual diffusions result in twice the complexity of INCRES, and moreover, the second PCUT diffusion depends on the first; they cannot be performed in parallel. While performing GROW in this way leads to a rigorous strategy for optimizing the likelihood $\mathcal{L}(\mathcal{P})$, in practice the second diffusion proves unnecessary. At the outset of the algorithm $\log \hat{U}$ dominates \tilde{U} by an order of magnitude, and approximating the PCUT algorithm by neglecting this term essentially results in INCRES. Making this approximation allows us to perform label propagation in a more efficient way, and as a consequence, we obtain an algorithm that runs more than twice as fast as the PCUT algorithm. The INCRES algorithm has a heuristic dynamic process rather than a rigorous energetic framework as its foundation, but as we show in the experimental section, the two approaches achieve comparable results in terms of accuracy.

The NIBBLE algorithm and its relatives [5, 9, 8, 7] also relate to INCRES in the sense that they obtain unsupervised clusterings from label propagation by planting random seeds. These works cluster the entire graph in a sequential manner: at each step a single random vertex is drawn and propagated. Then a sweep is performed to extract a small cluster around this vertex. These algorithms function well for problems aimed at extracting many small clusters from graphs with fine structure. In contrast, we perform multiway partitioning directly instead of recursively. We also aim at medium scale clusters instead of small scale clusters. We also utilize a significantly different random seeding strategy.

The INCRES algorithm alternates between label propagation (GROW) and thresholding (HARVEST). The idea of iteratively alternating between a few steps of label propagation and subsequent thresholding has also appeared in a transductive learning context [10], although the presence of labeled information results in a different implementation of the propagation step. The non-negative matrix factorization method [3] also incorporates random walk information in a manner that resembles the GROW routine, but otherwise the underlying principles of the algorithms differ substantially.

Finally, the algorithms GRACCLUS [1] and METIS [2] directly inspired the multigrid version of our algorithm. We use essentially the same coarsening algorithm, but rely upon a different clustering on the coarsest scale (INCRES vs. kernelized k -means or pure spectral clustering). Our refinement technique also differs substantially. The INCRES algorithm relates to the kernelized k -means procedure used in GRACCLUS even in the single level case: we can essentially interpret the GROW routine as the “maximization” step in an alternating minimization for a kernelized k -means. However, the kernel is a power of the normalized weights and the power may depend on the cluster, so it is not exactly the same. The “expectation” step in our algorithm is replaced by sampling, and instead of having a single representative for a class, the number of representatives increases as the algorithm progresses. Using power iterations of the weight matrix W directly for clustering has appeared in [11, 12]. These works utilize the power iterations to generate an em-

bedding of the vertices of the graph, which is then clustered using k -means. These methods can also be considered as kernelized k -means methods, with a power of the weights providing the kernel.

Because the GROW function we use iterates the random walk on the graph, our algorithm is a form of spectral clustering. However, our main contribution to the clustering problem, and the primary novelty in our algorithm, is the *incremental reseeding process*. This process is not fundamentally tied to the INGRES algorithm presented here— it seems to be quite universal and can be adapted to other clustering methods. However, combining reseeding with the random walk method offers an excellent combination of accuracy, speed, and robustness.

2.3 A Multigrid Speedup

The main computational burden of the basic INGRES algorithm lies in the GROW routine. Its computational cost scales like $O(R \times E \times \text{diam}(G))$ in the worst case, where E denotes the number of edges in the graph and $\text{diam}(G)$ denotes its diameter. In practice, the *expected* diameter of the graph effectively determines the cost of each step in the algorithm. For graphs commonly used in machine learning, such as k -nearest neighbor graphs, the number of matrix multiplications required in each call to GROW is generally quite small as a result.

However, the computational burden of the GROW routine still causes the straightforward implementation of our algorithm (in serial) to run two orders of magnitude slower than popular multiscale coarsen-and-refine algorithms such as GRACLUS [1] and METIS [2]. As we now show, pursuing a similar coarsen-and-refine strategy allows us remove an additional order of magnitude from the runtime of our algorithm. In many cases, this multiscale version of INGRES still maintains the consistently high level of accuracy obtained by the basic version.

Coarsening Phase: We follow the same coarsening procedure used by GRACLUS and its relatives [1, 2] in our multilevel approach. We construct an agglomerative hierarchy of weighted graphs in a recursive fashion, beginning from the original weighted graph. We therefore set $G^1 := G = (V, W)$ and then successively transform the vertex set $V^1 = V$ into a sequence of smaller weighted graphs G^2, G^3, \dots, G^L in such a way that the size of each corresponding vertex set decreases $|V^1| > |V^2| > \dots > |V^L|$ in a geometric fashion. The procedure terminates once the number of vertices $|V^L|$ in the current graph falls below some number n_0 , which we take as $n_0 = 20R$ in our experiments.

The transition between two successive levels G^l and G^{l+1} proceeds as follows. Each vertex of G^l begins unmarked. We then visit the vertices in V^l one-by-one according to their degree, from smallest to largest. When visiting a given vertex v_i , we merge it with its neighboring unmarked vertex v_j that maximizes

$$\frac{W_{ij}}{d_i} + \frac{W_{ij}}{d_j}.$$

The two merged vertices v_i and v_j are then marked, and the process continues. If at any stage a vertex v_i has no unmarked neighbors, we simply mark v_i and leave it as a singleton. The vertices v_i^{l+1} in G^{l+1} therefore correspond to pairs of vertices $\{v_{i_1}^l, v_{i_2}^l\}$ from G^l , so that the number of vertices $|V^l|$ roughly halves at each stage. Given two new vertices $v_i^{l+1} = \{v_{i_1}^l, v_{i_2}^l\}$ and $v_j^{l+1} = \{v_{j_1}^l, v_{j_2}^l\}$, we define the weights $W_{i,j}^{(l+1)}$ between these two vertices according to the relations

$$\begin{aligned} W_{ij}^{(l+1)} &:= W_{i_1 j_1}^{(l)} + W_{i_1 j_2}^{(l)} + W_{i_2 j_1}^{(l)} + W_{i_2 j_2}^{(l)}, \\ W_{ii}^{(l+1)} &:= W_{i_1 i_1}^{(l)} + 2W_{i_1 i_2}^{(l)} + W_{i_2 i_2}^{(l)}, \end{aligned}$$

i.e. simply by summing the weights between all possible pairs of vertices in the two merged pairs.

Base Clustering and Refinement: The clustering phase begins with a base clustering of the coarsest graph G^L in the hierarchy. We simply use the basic version of INGRES applied to G^L to obtain this initial clustering. We then extrapolate this clustering to the next level G^{L-1} in the hierarchy, refine the clustering of G^{L-1} using INGRES again, and repeat until we obtain a clustering of the original graph at the finest level.

The extrapolation procedure we use is straightforward: a clustering of G^{l+1} defines a corresponding clustering of G^l by simply assigning each vertex v_i^l and v_j^l in the pair $v_i^{l+1} = \{v_{i_1}^l, v_{i_2}^l\}$ to the cluster of its parent. We subsequently

Table 1 Algorithmic Comparison via Cluster Purity.

Data	size	R	RND	NCUT	LSD	NMFR	MTV	PCUT (speed 1)	INCRES (speed 1)	INCRES (speed 5)
20NEWS	20K	20	6%	27%	34%	61%	36%	61%	61%	61%
RCV1	9.6K	4	30%	38%	38%	43%	43%	53%	55%	51%
WEBKB4	4.2K	4	39%	40%	46%	58%	45%	58%	57%	57%
CITeseer	3.3K	6	22%	23%	53%	63%	43%	63%	62%	62%
MNIST	70K	10	11%	77%	76%	97%	96%	97%	96%	94%
PENDIGIT	11K	10	12%	80%	86%	87%	87%	87%	89%	86%
USPS	9.3K	10	17%	72%	70%	86%	85%	89%	88%	87%
OPTDIGIT	5.6K	10	12%	91%	91%	98%	95%	98%	97%	95%

refine this clustering using a slightly modified version the original INCRES procedure. Let $N^l := |V^l|$ denote the number of vertices in the current graph. We set the initial number of seeds to $s_i = .2N^l$, the final number of seeds to $s_f = .5N^l$, and we specify a set number of iterations I^l of INCRES to perform at the current level. We then select the seed increment parameter ds so that the number of seeds s transitions from s_i to s_f in exactly I^l iterations. We perform $I^L = 200$ steps of INCRES at the coarsest level and $I^1 = 1$ step of INCRES at the finest level. We then select the number of iterations I^l to perform at the l^{th} level as a geometrically decreasing progression between these two endpoints.

3 Experiments

We now provide the results of our extensive experimental evaluation of the algorithm. This section shows that our algorithm achieves state-of-the-art performance in terms of cluster purity on a variety of real word data sets while running faster than the other comparably accurate clustering methods. We also show that INCRES is very robust to perturbations in the input graph.

The Algorithms: We compare our method against five clustering algorithms that rely on variety of different principles. We select algorithms that, like our algorithm, partition the graph in a direct, non-recursive manner. The PCut algorithm [4] shares many of the features and motivations of the INCRES algorithm, but has twice the complexity per step. The NCut algorithm [13] is a widely used spectral algorithm that relies on a post-processing of the eigenvectors of the graph Laplacian to optimize the normalized cut energy. The NMFR algorithm [3] uses non-negative matrix factorization and graph-based random walk principles in order to factorize and regularize the original input similarity matrix. The LSD algorithm [14] provides another non-negative matrix factorization algorithm. It aims at finding a left-stochastic decomposition of the similarity matrix. The MTV algorithm from [15] provides a total-variation based algorithm that attempts to find an optimal multiway Cheeger cut of the graph by using ℓ^1 optimization techniques. The last three algorithms (NMFR, LSD and MTV) all use NCut in order to obtain an initial partition. By contrast, we initialize our algorithm with a random partition. We use the code available from [13] for NCut, the code available from [3] to test the two non-negative matrix factorization algorithms (NMFR and LSD) and the code available from [15] for the MTV algorithm.

The Data Sets: We provide experimental results on four text data sets (20NEWS, RCV1, WEBKB4, CITESEER) and four data sets containing images of handwritten digits (MNIST, PENDIGITS, USPS, OPTDIGITS). We processed the text data sets by removing a list of stop words as well as by removing all words with fewer than twenty occurrences (for 20NEWS) and fewer than five occurrences (for all others) across the corpus. We then construct a 5-NN graph based on the cosine similarity between tf-idf features. For variety, we include some weighted graphs (RCV1 and CITESEER) as well as some unweighted graphs (20NEWS and WEBKB4). For MNIST, PENDIGITS and OPTDIGITS we use the similarity matrices constructed by [3], where the authors first extract scattering features [16] for images before calculating an unweighted 10-NN graph. For USPS we constructed a weighted 10-NN graph from the raw data without any preprocessing. We provide the source for these data sets and more details on their construction in the supplementary material.

Accuracy Comparisons: In Table 1 we report the accuracy obtained by the selected algorithms NCUT LSD, NMFR, MTV, PCUT and INCRES (for two values of the timestep parameter, **speed** = 1 and **speed** = 5) on the various data sets. We use cluster purity to quantify the quality of the calculated partition, defined according to the relation

$$\text{Purity} = \frac{\text{number of "successes"}}{N} = \frac{1}{N} \sum_{r=1}^R \max_{1 \leq i < R} n_{r,i}.$$

Here $n_{r,i}$ denotes the number of data points in the r^{th} cluster that belong to the i^{th} ground-truth class. In other words, given a computed cluster we count a data point as a success if it belongs to the ground truth class that best represents the cluster. We allowed each iterative algorithm a total of 10,000 iterations to reach convergence. Both INCRES, PCUT and MTV rely on randomization, so for these algorithm we report the average purity achieved over 1000 different runs. The fourth column of the table (RND) provides a base-line purity for reference, i.e. the purity obtained by assigning each data point to a class from 1 to R uniformly at random. The boldface numbers in the table indicate the highest purity score achieved on each data set.

Overall, INCRES, PCUT and NMFR significantly outperform the other algorithms. This is especially true for text data sets. These three algorithms utilize a random walk strategy to help “smooth” irregular graphs, such as the similarity matrices obtained from text data sets. This strategy also contributes to the robustness of these algorithms and to their solid performance across the full range of data sets. However, the INCRES algorithm typically runs at least one order of magnitude faster than the NMFR algorithm and more than twice as fast as the PCUT algorithm. As that INCRES and PCUT obtain very similar purity scores, these results provide evidence of the fact that the sophisticated GROW routine of the PCUT algorithm does not lead to substantially better results than the simple and more efficient GROW routine of the INCRES algorithm. The additional mathematical rigor of PCUT does not translate into better results in practice, and it comes with a non-trivial increase in computational cost.

Finally, note that the INCRES algorithm performs comparably when **speed** = 1 and **speed** = 5, demonstrating that the algorithm is robust with respect to the choice of the seed increment parameter ds .

Speed Comparisons: Figure 2 illustrates the speed at which, LSD, MTV, NMFR and INCRES converge toward their respective solutions. We ran each algorithm for a total of 7 minutes on 20NEWS and for 15 minutes on MNIST. We report the purity obtained by the algorithm at each iteration. For the randomized algorithm (INCRES and MTV) the purity curves were obtained by averaging the results over 240 runs. The overwhelming computational burden for all of these algorithms arises from the sparse-matrix times full-matrix multiplications required at each step. The PCUT algorithm (not shown) runs about two times slower than the INCRES algorithm but otherwise achieves similar results. Each algorithm is implemented in a fair and consistent way, and the experiments were all performed on the same architecture.

Table 2 Computational Time

Data	NMFR	MTV	INC. (<i>spd1</i>)	INC. (<i>spd5</i>)
20NEWS	3.7mn (57.7%)	–	25.4s (57.7%)	5.6s (58%)
MNIST	4.6mn (92.2%)	1.8mn (90.7%)	4.8s (91.2%)	3.1s (89.3%)

In order to give an indication of the speed/accuracy trade-off for each algorithm, in Table 2 we record the time it took for the purity obtained by each algorithm to reach 95% of its limiting value on both 20NEWS and on MNIST. Overall, the simple INCRES algorithm provides accuracy comparable to the state-of-the-art NMF algorithm [3], yet runs an order of magnitude faster. Timing results on the data sets from table 1 are consistent with those obtained for 20NEWS and MNIST, in the sense that INCRES typically runs one order of magnitude faster than NMFR on these data sets as well.

Robustness Experiments: Table 3 reports accuracy results of various algorithms on graphs that we corrupted by adding different levels of noise. We began with the original 20NEWS graph used in Table 1 and added additional

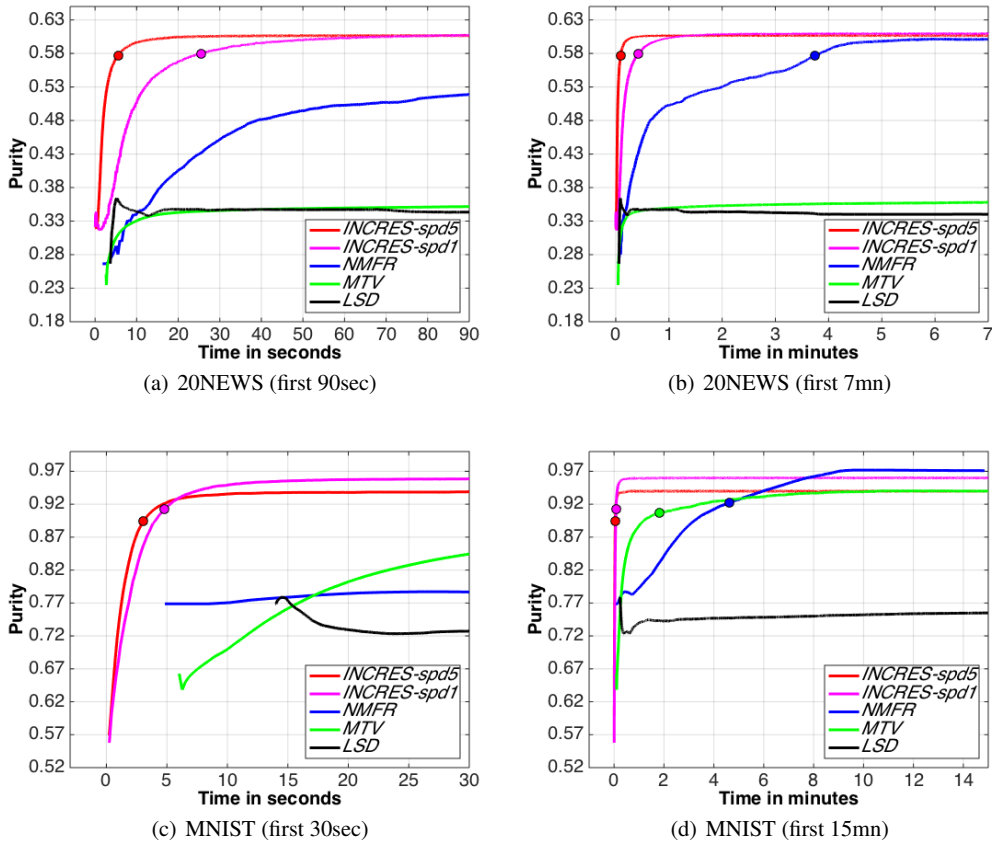


Fig. 2 Purity curves for the four algorithms considered on two benchmark data sets (20NEWS and MNIST). We plot purity against time for each algorithm over two different time windows. The circular marks on each curve indicate the point at which the curve reaches 95% of its limiting value. The corresponding times at which this happens are reported in Table 2.

edges to the graph uniformly at random. The original graph had $e = 144,632$ edges. For the experiment, we added $0.5e$, e , $1.5e$ and $2e$ additional noise edges. For each of these four levels of noise, we randomly generated 144 separate perturbed graphs. The table reports, for each level of noise, the average purity obtained by each algorithm on the 144 randomly generated matrices. We then proceeded to perturb the original MNIST graph in a similar fashion. The original graph has $e = 1,027,412$ edges, and we randomly generated 120 graphs at each level of noise. This gives a total of 1056 randomly generated graphs for this set of experiments. We provide experimental results for all algorithms other than NMFR (it is far too slow to run to convergence on all 1056 adjacency matrices) and PCUT (to avoid redundancy).

The results clearly elucidate the robustness of the INCREs algorithm with respect to noise in the graph construction process. On the 20NEWS data set, for example, all other algorithms experience a sharp decrease in accuracy as soon as noise is added. In contrast, the purity of the INCREs algorithm slowly decreases in a stable fashion. On the MNIST data sets, the results obtained by INCREs remain essentially unchanged across all noise levels. The competing algorithms do not exhibit this behavior. Interestingly, NCut and LSD actually obtain *better* results at the 50% and 100% noise levels. Given that LSD relies on NCut for initialization, it comes as no surprise that gains for NCut produce subsequent gains for LSD as well. This pathological behavior still indicates a lack of robustness, in the sense that both algorithms exhibit a high degree of sensitivity to changes in the underlying graph.

Multigrid Experiments: Table 4 reports the accuracies and run times of the coarsen and refine algorithms METIS [2], GRACLUS [1] and the multilevel version of our reseeding algorithm. We refer to the multilevel version of INCREs as INCREs-ml to distinguish it from the basic version. The reported purity values correspond to the average accuracy

Table 3 Robustness Comparisons

Noise	NCUT	LSD	MTV	GRACCLUS (multilevel)	INCREs-ml (multilevel)	INCREs (speed1)
20NEWS						
+0% edges	27%	34%	36%	42%	59%	61%
+50% edges	21%	27%	20%	15%	19%	52%
+100% edges	18%	22%	11%	13%	12%	44%
+150% edges	15%	20%	10%	12%	10%	34%
+200% edges	14%	18%	9%	11%	9%	27%
MNIST						
+0% edges	77%	76%	96%	97%	97%	96%
+50% edges	87%	94%	55%	93%	97%	97%
+100% edges	84%	93%	25%	80%	90%	97%
+150% edges	74%	87%	18%	63%	74%	97%
+200% edges	67%	82%	16%	52%	53%	96%

Table 4 Accuracy Comparison for Multilevel Algorithms.

Data	size	METIS	GRACCLUS	INCREs-ml
20NEWS	20K	42.4%	42.4% (0.05s)	58.6% (0.85s/0.07s)
RCV1	9.6K	34.1%	42.4% (0.01s)	47.9% (0.15s/0.03s)
WEBKB4	4.2K	37.9%	49.0% (0.01s)	53.9% (0.10s/0.02s)
CITeseer	3.3K	45.2%	53.5% (0.01s)	61.5% (0.11s/0.02s)
MNIST	70K	86.0%	96.9% (0.17s)	96.9% (1.99s/0.52s)
PENdIGIT	11K	67.3%	84.7% (0.02s)	88.8% (0.46s/0.05s)
USPS	9.3K	75.1%	86.9% (0.02s)	87.2% (0.34s/0.05s)
OPTIDIGIT	5.6K	83.0%	94.2% (0.01s)	95.0% (0.26s/0.03s)

obtained over 500 trials of each routine. By and large, INCREs-ml obtains clustering of higher quality than the two other multilevel algorithms. This comes at a cost of one order of magnitude in computational time.

The computational cost that INCREs-ml incurs at level l of the graph hierarchy scales as

$$O(R \times E_l \times \text{diam}(G_l) \times I^l),$$

where E_l denotes the number of edges in the graph G^l and I^l once again denotes the number of iterations performed at this level. In comparison, the computational burden that GRACCLUS incurs at level l scales as

$$O(R \times E_l \times I^l),$$

where I^l has the same meaning as before, but obviously refers to the number of iterations performed by GRACCLUS. The extra term $\text{diam}(G_l)$ in the computational cost of INCREs-ml partly explains the difference between the computational times. However, the diameters of the graphs considered in our experiments are typically smaller than ten. This is especially true for the coarsest levels in the hierarchy. The other major source of the difference in computational time comes from implementation: while GRACCLUS is a heavily optimized C code, our multilevel algorithm has a naive MATLAB implementation. Finally, we remark that all of these multilevel procedures suffer from a severe sensitivity to noise, see Table 3. This fact motivates the use of the basic version of INCREs in situations where robustness to the graph construction process is highly desirable.

4 Appendix

4.1 Matlab Code used in the experimental section

Below is the exact MATLAB code that we used in the experimental section of the paper.

```
1 function C=INCREAS( W , R , speed , maxiter)
2 D = sum(W,2);
3 RW=W*spdiags(1./D,0,n,n);
4 s=1;
5 ds=speed*10^(-4)*n/R;
6 C=randi(R,n,1);
7 for iter=1:maxiter
8     [F,R] = PLANT(C, R, round(s));
9     while ( min( min(F) ) < eps )
10        F= RW * F;
11    end
12    [~,C] = max(F, [],2);
13    s=s+ds;
14 end
15 end
```

```
1 function [F,R] = PLANT(C,R,s)
2 n=length(C);
3 F=zeros(n,R);
4 EmptyClass=[];
5 for k=1:R
6     idx= find(C==k);
7     ClassSize=length(idx);
8     if (ClassSize== 0)
9         EmptyClass=[EmptyClass,k];
10    else
11        idxSeeds = idx( randi(ClassSize,s,1) );
12        F(:,k)= full(sparse(idxSeeds,ones(1,s),1,n,1));
13    end
14 end
15 F(:,EmptyClass)=[];
16 R=size(F,2);
17 end
```

4.2 Datasets

- 20NEWS (unweighted similarity matrix): The word count from the raw documents was computed using the Rainbow library [17] with a default list of stop words. Words appearing less than 20 times were also removed. The similarity matrix was then obtained by 5 nearest neighbors using cosine similarity between tf-idf features. Source: <http://www.cs.cmu.edu/~mccallum/bow/rainbow/>
- RCV1 (weighted similarity matrix): This dataset was obtained in preprocessed format from <http://www.cad.zju.edu.cn/home/dengcai/Data/TextData.html> with the tf-idf features were already computed. We then simply used cosine similarity and 5-NN.
- WEBKB4 (unweighted similarity matrix): The word count from the raw documents was done with the Rainbow library [17]. A list of stop word was removed. Words appearing less than 5 times were removed. The similarity

matrix was then obtained by 5 nearest neighbors using cosine similarity between tf-idf features. Source: <http://www.cs.cmu.edu/afs/cs/project/theo-20/www/data/>

- CITESEER (weighted similarity matrix): This dataset was obtained in preprocessed format from <http://lings.cs.umd.edu/projects//projects/lbc/index.html> where each publication in the dataset is described by a 0/1-valued word vector indicating the absence/presence of the corresponding word from the dictionary. We then simply used cosine similarity and 5-NN.
- MNIST, PENDIGITS, OPTDIGITS (unweighted similarity matrix): The similarity matrices were obtained from [3], where the authors first extracted scattering features using [18] for images before calculating the 10-NN graph. Source: <http://users.ics.aalto.fi/rozyang/nmfr/index.shtml>
- USPS (weighted similarity matrix): We computed a 10-NN graph using standard Euclidean distance between the raw images. Each edge in the 10-NN graph was given the weight

$$w_{ij} = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

where each x_i denotes a vector containing the raw pixel data. The parameter σ was chosen as the mean distance between each vertex and its 10th nearest neighbor. Source: <http://www.cad.zju.edu.cn/home/dengcai/Data/MLData.html>

Acknowledgement

XB is supported by NRF Fellowship NRFF2017-10.

References

1. Inderjit S. Dhillon, Yuqiang Guan, and Brian Kulis. Weighted graph cuts without eigenvectors: A multilevel approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(11):1944–1957, 2007.
2. George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.*, 20(1):359–392, December 1998.
3. Zhirong Yang, Tele Hao, Onur Dikmen, Xi Chen, and Erkki Oja. Clustering by nonnegative matrix factorization using graph random walk. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1088–1096, 2012.
4. Xavier Bresson, Thomas Laurent, Arthur Szlam, and James H von Brecht. The product cut. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
5. László Lovász and Miklós Simonovits. Random walks in a convex body and an improved volume algorithm. *Random structures & algorithms*, 4(4):359–412, 1993.
6. Xiaojin Zhu, Zoubin Ghahramani, and John Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *IN ICML*, pages 912–919, 2003.
7. Daniel A. Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and its application to nearly linear time graph partitioning. *SIAM Journal on Computing*, 42(1):1–26, 2013.
8. Reid Andersen, Fan Chung, and Kevin Lang. Local graph partitioning using pagerank vectors. In *Proceedings of the 47th Annual Symposium on Foundations of Computer Science (FOCS '06)*, pages 475–486, 2006.
9. Daniel A. Spielman and Shang-Hua Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 81–90, 2004.
10. Cristina Garcia-Cardona, Ekaterina Merkurjev, Andrea L. Bertozzi, Arjuna Flenner, and Allon G. Percus. Multiclass data segmentation using diffuse interface methods on graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 99:1, 2014.
11. Frank Lin and William W. Cohen. Power iteration clustering. In *ICML*, pages 655–662, 2010.
12. Stéphane Lafon and Ann B. Lee. Diffusion maps and coarse-graining: A unified framework for dimensionality reduction, graph partitioning, and data set parameterization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 28(9):1393–1403, 2006.
13. Stella X. Yu and Jianbo Shi. Multiclass spectral clustering. in international conference on computer vision. In *International Conference on Computer Vision*, 2003.
14. Raman Arora, M Gupta, Amol Kapila, and Maryam Fazel. Clustering by left-stochastic matrix factorization. In *International Conference on Machine Learning (ICML)*, pages 761–768, 2011.
15. Xavier Bresson, Thomas Laurent, David Uminsky, and James von Brecht. Multiclass total variation clustering. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.

16. Joan Bruna and Stéphane Mallat. Invariant scattering convolution networks. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(8):1872–1886, 2013.
17. Andrew Kachites McCallum. Bow: A toolkit for statistical language modeling, text retrieval, classification and clustering. <http://www.cs.cmu.edu/mccallum/bow>, 1996.
18. M. Stephane. Group invariant scattering. *Communications on Pure and Applied Mathematics*, 65(10):1331–1398, 2012.