Graph Degree Linkage: Agglomerative Clustering on a Directed Graph

Wei Zhang¹, Xiaogang Wang^{2,3}, Deli Zhao¹, and Xiaoou Tang^{1,3}

- Department of Information Engineering, The Chinese University of Hong Kong wzhang009@gmail.com
- ² Department of Electronic Engineering, The Chinese University of Hong Kong
- ³ Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, China

Abstract. This paper proposes a simple but effective graph-based agglomerative algorithm, for clustering high-dimensional data. We explore the different roles of two fundamental concepts in graph theory, indegree and outdegree, in the context of clustering. The average indegree reflects the density near a sample, and the average outdegree characterizes the local geometry around a sample. Based on such insights, we define the affinity measure of clusters via the product of average indegree and average outdegree. The product-based affinity makes our algorithm robust to noise. The algorithm has three main advantages: good performance, easy implementation, and high computational efficiency. We test the algorithm on two fundamental computer vision problems: image clustering and object matching. Extensive experiments demonstrate that it outperforms the state-of-the-arts in both applications. I

Introduction

Many problems in computer vision involve clustering. Partitional clustering, such as k-means [1], determines all clusters at once, while agglomerative clustering [1] begins with a large number of small clusters, and iteratively selects two clusters with the largest affinity under some measures to merge, until some stopping condition is reached. Agglomerative clustering has been studied for more than half a century, and used in many applications [1], because it is conceptually simple and produces an informative hierarchical structure of clusters.

Classical agglomerative clustering algorithms have several limitations [1], which have restricted their wider applications in computer vision. The data in computer vision applications are usually high dimensional. The distributions of data clusters are often in different densities, sizes, and shapes, and form manifold structures. In addition, there are often noise and outliers in data. The conventional agglomerative clustering algorithms, such as the well-known linkage methods [1], usually fail to tackle these challenges. As their affinities are directly computed using pairwise distances between samples and cannot capture the global manifold structures in high-dimensional spaces,

¹ The code and supplemental materials are publicly available at http://mmlab.ie.cuhk.edu.hk/research/gdl/.

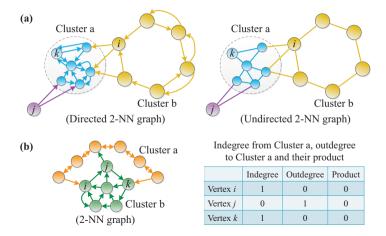


Fig. 1. (a) Indegree can be use to detect the change of densities. The density in Cluster a is high, and the density in Cluster b is low. The vertices inside Cluster a are strongly connected, but there is no outedge to vertices outside Cluster a. So the indegree of k from Cluster a is nonzero, while the indegree of i (a vertex in Cluster b) and j (an outlier) from Cluster a are zero. If an undirected graph is considered without separating indegrees and outdegrees, both i and j have the same degree from Cluster a as k. (b) The product of the indegree and outdegree is an affinity measure robust to noisy edges between the two clusters. Under this measure, Cluster a and Cluster b have a zero affinity, i.e., the sum of product of indegree and outdegree for all vertices is 0, and thus they are separated well.

these algorithms have problems of clustering high-dimensional data, and are quite sensitive to noise and outliers [1].

To tackle these problems, we propose a simple and fast graph-based agglomerative clustering algorithm. The graph representation of data has been extensively exploited in various machine learning topics [2,3,4,5,6], but has rarely been utilized in agglomerative clustering. Our algorithm builds K-nearest-neighbor (K-NN) graphs using the pairwise distances between samples, since studies [4] show the effectiveness of using local neighborhood graphs to model data lying on a low-dimensional manifold embedded in a high-dimensional space.

We use the *indegree* and *outdegree*, fundamental concepts in graph theory, to characterize the affinity between two clusters. The outdegree of a vertex to a cluster measures the similarity between the vertex and the cluster. If many of the K-NNs of the vertex belong the cluster, the outdegree is large. The outdegree can capture the manifold structures in the high dimensional space. The indegree of a vertex from a cluster reflects the density near the vertex. It is effective for detecting the change of densities, which often occurs at the boundary of clusters. Therefore, we use it to separate clusters close in space but different in densities, and also reduce the effect of noise. An example is shown in Fig. 1(a). To our best knowledge, properties of the indegree and outdegree have not been explored by any existing clustering algorithm, although they were successfully applied in analysis of complex networks such as World Wide Web [9] and social networks [10] and showed interesting results.

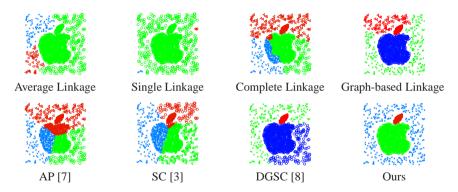


Fig. 2. Results of different clustering algorithms on a synthetic multiscale dataset. Our algorithm can perfectly discover the three clusters with different shapes, sizes, and densities. The output clusters are shown in color (best viewed on screen).

Our affinity measure between two clusters is defined as follows. First, the structural affinity from a vertex to a cluster is defined via the product of the average indegree from the cluster and average outdegree to the cluster. Intuitively, if a vertex belongs to a cluster, it should be strongly connected to the cluster, i.e., both its indegree and outdegree are large. Otherwise, either the indegree or outdegree is small. Therefore, the product of indegree and outdegree can be a good affinity measure (Fig. 1(b)). We show that the correlation between the inter-cluster indegree and outdegree is weak across different vertices, if the two clusters belong to different ground-truth clusters, using synthetic data in Fig. 3. Then, the affinity between two clusters is naturally the aggregated affinity measure for all the vertices in the two clusters.

Our algorithm has three main advantages as follows.

First of all, it has *outstanding performance*, especially on noisy data and multiscale data (i.e., clusters in different densities). The visual comparisons with linkage methods [1], graph-based average linkage, affinity propagation (AP) [7], spectral clustering (SC) [3], and directed graph spectral clustering (DGSC) [8] on synthetic multiscale data are shown in Fig. 2. Noise and multiple scales can degrade the performance of spectral clustering greatly [11], while the indegree and outdegree in our algorithm detect the boundary of scales automatically² and reduce the effect of noise. In Sec. 4, extensive experiments on real data, including imagery data and feature correspondence data, demonstrate its superiority over state-of-the-art methods. These experiments aim at two fundamental problems in computer vision, i.e., image clustering and object matching, and the results suggest many potential applications of our work.

Second, it is *easy to implement*. This affinity measure can be expressed in a matrix form and implemented with vector additions and inner-products. Therefore, our algorithm can be implemented without any dependency on external numerical libraries,

 $^{^2}$ E.g., if cluster a has higher density than cluster b, the boundary of cluster a will have high indegree and low outdegree, while the boundary of cluster b will have low indegree and high outdegree.

such as eigen-decomposition which was extensively employed by many clustering algorithms [2,3,12].

Finally, it is *very fast*. We propose an acceleration method for our algorithm. In practice, our algorithm is much faster than spectral clustering [2,3], especially on large-scale data.

2 Related Work

The literature dedicated to agglomerative clustering is abundant [1,13,14]. Linkages [1], e.g., average linkage, define the affinity based on pairwise distances between samples. Since pairwise distances do not well capture the global structures of data, these methods fail on clustering data with complex structures and are sensitive to noise [1] (see the example in Fig. 2). Many variants of linkage methods, such as DBSCAN [15], have been proposed in the data mining community and show satisfactory performance. However, they usually fail to tackle the great challenge from high-dimensional spaces, because their sophisticated affinity measures are based on observations from low-dimensional data [16].

Several algorithms [17,18,19] has attempted to perform agglomerative clustering on the graph representation of data. Chameleon [17] defines the cluster affinity from relative interconnectivity and relative closeness, both of which are based on a min-cut bisection of clusters. Although good performance was shown on 2D toy datasets, it suffers from high computational cost because its affinity measure is based on a min-cut algorithm. Zell [18] describes the structure of a cluster via the zeta function and defined the affinity based on the structural changes after merging. It needs to compute matrix inverse in each affinity computation, so it is much slower than our simple algorithm (see Sec. 4.1). Felzenszwalb and Huttenlocher proposed an effective algorithm for image segmentation [19].

Besides agglomerative clustering, K-means [1] and spectral clustering [2,3,20] are among the most widely used clustering algorithms. However, K-means is sensitive to the initialization and difficult to handle clusters with varying densities and sizes, or manifold shapes. Although spectral clustering can handle the manifold data well, its performance usually degrades greatly with the existence of noise and outliers, because the eigenvectors of graph Laplacian are sensitive to noisy perturbations [5]. Affinity Propagation [7] explores the intrinsic data structures by message passing among data points. Although it performs well on high-dimensional data, it usually requires considerable run-time, especially when the preference value cannot be manually set.

Directed graphs have been studied for spectral clustering (e.g., [8]). However, these methods symmetrize the directed graph before the clustering task. In contrast, we only symmetrize the affinity between two clusters, while keep the directed graph during the clustering process. Therefore, our algorithm utilizes more information from the asymmetry and is more robust to noisy edges (see Fig. 2 for a comparison between DGSC [8] and our algorithm).

3 Graph Degree Linkage

3.1 Neighborhood Graph

Given a set of samples $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$, we build a directed graph G = (V, E), where V is the set of vertices corresponding to the samples in \mathcal{X} , and E is the set of edges connecting vertices. The graph is associated with a weighted adjacency matrix $\mathbf{W} = [w_{ij}]$, where w_{ij} is the weight of the edge from vertex i to vertex j. $w_{ij} = 0$ if and only if there is no edge from i to j.

To capture the manifold structures in high-dimensional spaces, we use the K-NN graph, in which the weights are defined as

$$w_{ij} = \begin{cases} \exp\left(-\frac{dist(i,j)^2}{\sigma^2}\right), & \text{if } \mathbf{x}_j \in \mathcal{N}_i^K, \\ 0, & \text{otherwise,} \end{cases}$$
 (1)

where dist(i,j) is the distance between \mathbf{x}_i and \mathbf{x}_j , \mathcal{N}_i^K is the set of K-nearest neighbors of \mathbf{x}_i , and σ^2 is set as $\sigma^2 = \frac{a}{nK} \left[\sum_{i=1}^n \sum_{\mathbf{x}_j \in \mathcal{N}_i^K} dist(i,j)^2 \right]$. K and a are free parameters to be set. In a K-NN graph, there is an edge pointing from \mathbf{x}_i to \mathbf{x}_j with weight w_{ij} , if $\mathbf{x}_j \in \mathcal{N}_i^K$.

3.2 Algorithm Overview

The graph degree linkage (GDL) algorithm begins with a number of initial small clusters, and iteratively selects two clusters with the maximum affinity to merge. The affinities are computed on the K-NN graph, based on the indegree and outdegree of vertices in the two clusters.

The initial small clusters are simply constructed as weakly connected components of a K^0 -NN graph, where the neighborhood size K^0 is small, typically as 1 or 2. Then, each component is an initial cluster, and each sample is assigned to only one cluster.

Definition 1 A connected component of an undirected graph is a maximal connected subgraph in which any two vertices are connected to each other by paths.

A weakly connected component of a directed graph is a connected component of the undirected graph produced by replacing all of its directed edges with undirected edges.

The GDL algorithm is presented as Algorithm 1, with details given in the following subsection.

3.3 Affinity Measure via Product of Indegree and Outdegree

The affinity measure between two clusters is the key of an agglomerative clustering algorithm. Our affinity measure is based on *indegree* and *outdegree* in the graph representation. For simplicity, we start from measuring the affinity between a vertex and a cluster.

Indegree and outdegree. Considering a vertex and a cluster, the connectivity between them by inedges and outedges can be quantified using the concepts of indegree and outdegree.

Algorithm 1 Graph Degree Linkage (GDL)

Input: a set of n samples $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n\}$, and the target number of clusters n_T .

Build the K^0 -NN graph, and detect its weakly connected components as initial clusters. Denote the set of initial clusters as $V^c = \{C_1, \dots, C_{n_c}\}$, where n_c is the number of clusters.

Build the K-NN graph, and get the weighted adjacency matrix \mathbf{W} .

while $n_c > n_T$ do

Search two clusters C_a and C_b , such that $\{C_a, C_b\} = \operatorname{argmax}_{C_a, C_b \in V^c} \mathcal{A}_{C_a, C_b}$, where \mathcal{A}_{C_a, C_b} is the affinity measure between C_a and C_b , computed using Eq. (5).

 $V^c \leftarrow \{V^c \setminus \{\mathcal{C}_a, \mathcal{C}_b\}\} \cup \{\mathcal{C}_a \cup \mathcal{C}_b\}, \text{ and } n_c = n_c - 1.$

end while Output: V^c .

Definition 2 Given a vertex i, the average indegree from and the average outdegree to a cluster C is defined as $\deg_i^-(C) = \frac{1}{|C|} \sum_{j \in C} w_{ji}$ and $\deg_i^+(C) = \frac{1}{|C|} \sum_{j \in C} w_{ij}$, respectively, where |C| is the cardinality of set C.

As we stated in Sec. 1, the indegree measures the density near sample i, and the out-degree characterizes the K-NN similarity from vertex i to cluster \mathcal{C} . We use the size of the cluster to normalize the degrees, otherwise, the algorithm may favor of merging large clusters instead of merging small clusters with dense connections. We find that in practice the normalized degrees work much better than the unnormalized degrees.

Affinity between a vertex and a cluster. A vertex should be merged to a cluster if it is strongly connected to the cluster by both inedges and outedges. Mathematically, the correlation of two types of degree is weak, if the vertex and the cluster belong to different ground-truth clusters, and strong, otherwise. To verify this intuition, we show such statistics on synthetic data in Fig. 3. Therefore, we define the affinity as the product of the average indegree and average outdegree, i.e.,

$$\mathcal{A}_{i\to\mathcal{C}} = \deg_i^-(\mathcal{C}) \deg_i^+(\mathcal{C}). \tag{2}$$

This affinity is robust to noisy edges between different ground-truth clusters because the product can be zero if the inedges and outedges do not coincide.

Affinity between two clusters. Following the above, we define the asymmetric affinity from cluster C_b to cluster C_a by summing up with respect to all the vertices in C_b , i.e.,

$$\mathcal{A}_{\mathcal{C}_b \to \mathcal{C}_a} = \sum_{i \in \mathcal{C}_b} \mathcal{A}_{i \to \mathcal{C}_a} = \sum_{i \in \mathcal{C}_b} \deg_i^-(\mathcal{C}_a) \deg_i^+(\mathcal{C}_a). \tag{3}$$

Finally, we have the symmetric affinity used in our algorithm as

$$\mathcal{A}_{\mathcal{C}_a,\mathcal{C}_b} = \mathcal{A}_{\mathcal{C}_b \to \mathcal{C}_a} + \mathcal{A}_{\mathcal{C}_a \to \mathcal{C}_b} \tag{4}$$

Efficient computation of affinity. Our affinity measure can be computed efficiently using the following theorem.

Theorem 1 The affinity between C_a and C_b defined in Eq. (4) can be expressed in the matrix form

$$\mathcal{A}_{\mathcal{C}_a,\mathcal{C}_b} = \frac{1}{|\mathcal{C}_a|^2} \mathbf{1}_{|\mathcal{C}_a|}^T \mathbf{W}_{\mathcal{C}_a,\mathcal{C}_b} \mathbf{W}_{\mathcal{C}_b,\mathcal{C}_a} \mathbf{1}_{|\mathcal{C}_a|} + \frac{1}{|\mathcal{C}_b|^2} \mathbf{1}_{|\mathcal{C}_b|}^T \mathbf{W}_{\mathcal{C}_b,\mathcal{C}_a} \mathbf{W}_{\mathcal{C}_a,\mathcal{C}_b} \mathbf{1}_{|\mathcal{C}_b|}, \quad (5)$$

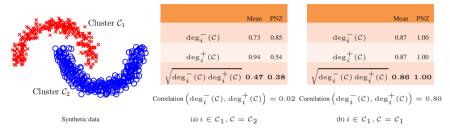


Fig. 3. To verify the robustness of the product of indegree and outdegree as an affinity measure from a vertex i to a cluster \mathcal{C} , we compare statistics in two cases: i and \mathcal{C} belong to different ground-truth clusters, e.g., $i \in \mathcal{C}_1$ and $\mathcal{C} = \mathcal{C}_2$ as in (a), and i and i and i are in the same ground-truth cluster, e.g., $i \in \mathcal{C}_1$ and i and i and i and i are in the same ground-truth cluster, e.g., $i \in \mathcal{C}_1$ and i and i and i and i are in the same ground-truth cluster, e.g., $i \in \mathcal{C}_1$ and i are in the same ground-truth cluster, e.g., $i \in \mathcal{C}_1$ and i are in the same ground-truth cluster, e.g., $i \in \mathcal{C}_1$ and i and i and i are in the same ground-truth clusters. For all $i \in \mathcal{C}_1$, such that i are much smaller than those of i and i and i and i are much smaller than those of i and i and i and i are much smaller than those of i and i and i and i are close to those of i and i and i and i and i and i and i are close to those of i and i

where $\mathbf{W}_{\mathcal{C}_a,\mathcal{C}_b}$ is the submatrix of \mathbf{W} whose row indices correspond to the vertices in \mathcal{C}_a and column indices correspond to the vertices in \mathcal{C}_b , i.e., the weights of edges from \mathcal{C}_a to \mathcal{C}_b , and $\mathbf{1}_L$ is an all-one vector of length L.

Remark 1 The computation is reduced to vector additions and inner-products. So, our algorithm is easy to implement.

Proof. It is easy to see that

$$\deg_i^-(\mathcal{C}_a) = \frac{1}{|\mathcal{C}_a|} \left[\mathbf{1}_{|\mathcal{C}_a|}^T \mathbf{W}_{\mathcal{C}_a, \mathcal{C}_b} \right]_i, \tag{6}$$

$$\deg_i^+(\mathcal{C}_a) = \frac{1}{|\mathcal{C}_a|} \left[\mathbf{W}_{\mathcal{C}_b, \mathcal{C}_a} \mathbf{1}_{|\mathcal{C}_a|} \right]_i, \tag{7}$$

where $[\mathbf{v}]_i$ is the *i*-th element of vector \mathbf{v} . Then, by Eq. (3), we can obtain the following lemma.

Lemma 2

$$\mathcal{A}_{\mathcal{C}_b \to \mathcal{C}_a} = \frac{1}{|\mathcal{C}_a|^2} \mathbf{1}_{|\mathcal{C}_a|}^T \mathbf{W}_{\mathcal{C}_a, \mathcal{C}_b} \mathbf{W}_{\mathcal{C}_b, \mathcal{C}_a} \mathbf{1}_{|\mathcal{C}_a|}.$$
 (8)

Finally, Theorem 1 can be directly implied by Lemma 2 using Eq. (4).

Comparison to average linkage. The GDL algorithm is different from average linkage in the following three aspects. First of all, the conventional average linkage is based on pairwise distances [1]. Although we find that average linkage has much better

performance on the K-NN graph than pairwise distances, we are unaware of any literature which studied the graph-based average linkage algorithm. Second, graph-based average linkage simply symmetrizes the directed graph by setting $w_{ij} = w_{ji} = (w_{ij} + w_{ji})/2$, while our algorithm uses the directed graph. Third, graph-based average linkage can be interpreted as defining the affinity measure $\mathcal{A}_{\mathcal{C}_b \to \mathcal{C}_a} = \frac{1}{|\mathcal{C}_b|} \sum_{i \in \mathcal{C}_b} [\deg_i^-(\mathcal{C}_a) + \deg_i^+(\mathcal{C}_a)]/2$ using our indegree-outdegree framework. The sum of the indegree and outdegree is not as robust as the product of them to noise. Experimental results in Fig. 2 and Sec. 4.1 demonstrate the superiority of GDL to graph-based average linkage.

3.4 Implementations of GDL

We present two implementations of the GDL algorithm: an exact algorithm via an efficient update formula and an approximate algorithm called Accelerated GDL (AGDL). Both implementations have the time complexity of $O(n^2)$ (see Theorem 3).

Update formula. In each iteration, we select two clusters C_a and C_b with the largest affinity and merge them as $C_{ab} = C_a \cup C_b$. Then, we need to update the asymmetric affinity $A_{C_{ab} \to C_c}$ and $A_{C_c \to C_{ab}}$, for any other cluster C_c .

Using Lemma 2, we find that $\mathcal{A}_{\mathcal{C}_{ab} \to \mathcal{C}_c}$ can be computed as follows.

$$\mathcal{A}_{\mathcal{C}_{ab} \to \mathcal{C}_c} = \mathcal{A}_{\mathcal{C}_a \to \mathcal{C}_c} + \mathcal{A}_{\mathcal{C}_b \to \mathcal{C}_c}. \tag{9}$$

By storing all the asymmetric affinities, the update is simple.

As the same update formula cannot be applied to $\mathcal{A}_{\mathcal{C}_c \to \mathcal{C}_{ab}}$, we have to compute it directly using Eq. (8). However, the total complexity is O(n) in each iteration, due to the row sparsity of \mathbf{W} (see Sec. 7 in the supplemental materials for details).

The GDL algorithm with the update formula (GDL-U) is presented as Algorithm 2 in the supplemental materials.

Accelerated GDL. Although the GDL-U algorithm is simple and fast, we further propose AGDL. The major computational cost is on computing the affinities. To reduce the number of affinities computed in each iteration, AGDL maintains a neighbor set of size K^c for each cluster in V^c , to approximate its K^c -nearest cluster set. Then, finding the maximum affinity among all pairs of clusters can then be approximated by searching it in all the neighbor sets. Updating the neighbor sets involves computation of the affinity between the new cluster and a small set of clusters, instead of all the other clusters.

Denote the neighbor set of a cluster \mathcal{C} as $\mathcal{N}_{\mathcal{C}}$. Initially $\mathcal{N}_{\mathcal{C}}$ consists of \mathcal{C} 's K^c -nearest clusters. Once two clusters \mathcal{C}_a and \mathcal{C}_b are merged, we need to update the neighbor sets which include \mathcal{C}_a or \mathcal{C}_b , and create the neighbor set of $\mathcal{C}_a \cup \mathcal{C}_b$. We utilize two assumptions that (1) if \mathcal{C}_a or \mathcal{C}_b is among the K^c -nearest clusters of \mathcal{C}_c ; (2) if \mathcal{C}_c is among the K^c -nearest clusters of \mathcal{C}_a or \mathcal{C}_b , \mathcal{C}_c is probably among the K^c -nearest clusters of $\mathcal{C}_a \cup \mathcal{C}_b$. So, the new cluster $\mathcal{C}_a \cup \mathcal{C}_b$ is added to the neighbor sets which include \mathcal{C}_a or \mathcal{C}_b previously. To create the neighbor set for $\mathcal{C}_a \cup \mathcal{C}_b$, we select the K^c -nearest clusters from $\mathcal{N}_{\mathcal{C}_a} \cup \mathcal{N}_{\mathcal{C}_b}$.

The AGDL algorithm is summarized in Algorithm 3 in the supplemental materials.

3.5 Time Complexity Analysis

We have the following theorem about the time complexity of the GDL, GDL-U and AGDL algorithms (please refer to Sec. 7 in the supplemental materials for the proof).

Theorem 3

- (a) The time complexity of the GDL algorithm (i.e., Algorithm 1) is $O(n^3)$.
- (b) The time complexity of the GDL-U algorithm (i.e., Algorithm 2) is $O(n^2)$.
- (c) The time complexity of the AGDL algorithm (i.e., Algorithm 3) is $O(n^2)$.

4 Experiments

In this section, we demonstrate the effectiveness of GDL and AGDL on image clustering and object matching. All the experiments are run in MATLAB on a PC with 3.20GHz CPU and 8G memory.

4.1 Image Clustering

We carry out experiments on six publicly available image benchmarks, including object image databases (COIL-20 and COIL-100), hand-written digit databases (MNIST and USPS), and facial image databases (Extended Yale-B, FRGC ver2.0).³ For MNIST, we use all the images in the testing set. For FRGC ver2.0, we use all the facial images in the training set of experiment 4. The statistics of all the datasets are presented in Table 2. We adopt widely used features for different kinds of images: the intensities of pixels as features and Euclidean distance for object and digit images, and local binary patterns (LBP) as features and χ^2 distance for facial images.

We compare the GDL-U and AGDL with eight representative algorithms, i.e., k-medoids (k-med) [1], average linkage (Link) [1], graph-based average linkage (G-Link), normalized cuts (NCuts) [2]⁴, NJW spectral clustering (NJW-SC) [3], directed graph spectral clustering (DGSC) [8], self-tuning spectral clustering (STSC) [11] and Zell [18]. Here we use k-medoids instead of k-means because it can handle the case where distances between points are not measured by Euclidean distances. To fairly compare the graph-based algorithms, we fix K=20 and select a with the best performance from the set $\{10^i, i \in [-2:0.5:2]\}$ on all the datasets. For our algorithms, the parameters are fixed as $K^0=1$, $K^c=10$. The numbers of ground-truth clusters are used as the input of all algorithms (e.g., n_T in our algorithm).

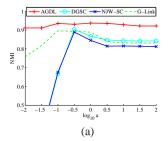
We adopt the widely used Normalized Mutual Information (NMI) [12] to quantitatively evaluate the performance of clustering algorithms. The NMI quantifies the normalized statistical information shared between two distributions. A larger NMI value indicates a better clustering result.

GOIL-20 and COIL-100 are from http://www.cs.columbia.edu/CAVE/software/.
MNIST and USPS are from http://www.cs.nyu.edu/~roweis/data.html. Extended Yale-B is from http://vision.ucsd.edu/~leekc/ExtYaleDatabase/ExtYaleBFRGC ver2.0 is from http://face.nist.gov/frgc/.

⁴ The code is downloaded from http://www.cis.upenn.edu/~jshi/software/, which implements the multiclass normalized cuts algorithm [20].

Table 1. Quantitative clustering results in NMI on real imagery data. A larger NMI value indicates a better clustering result. The results shown in a boldface are significantly better than the others, with a significance level of 0.01.

Dataset	k-med	Link	G-Link	NCuts	NJW-SC	DGSC	STSC	Zell	GDL-U	AGDL
COIL-20	0.710	0.647	0.896	0.884	0.889	0.904	0.895	0.911	0.937	0.937
COIL-100	0.706	0.606	0.855	0.823	0.854	0.858	0.858	0.913	0.929	0.933
USPS	0.336	0.095	0.732	0.675	0.690	0.747	0.726	0.799	0.824	0.824
MNIST	0.390	0.304	0.808	0.753	0.755	0.795	0.756	0.768	0.844	0.844
Yale-B	0.329	0.255	0.766	0.809	0.851	0.869	0.860	0.781	0.910	0.910
FRGC	0.541	0.570	0.669	0.720	0.723	0.732	0.729	0.653	0.747	0.746



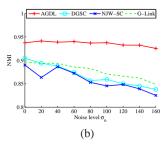


Fig. 4. Variations of performance of different clustering algorithms on the COIL-20 dataset, (a) when the parameter a for controlling σ in Eq. (1) changes; (b) when we add Gaussian noise $\mathcal{N}(0, \sigma_n^2)$ to the images. The NMI differences between $\sigma_n = 0$ and $\sigma_n = 160$ are 0.048, 0.065, 0.067, 0.012, for G-Link, NJW-SC, DGSC, and AGDL, respectively.

The results measured in NMI are given in Table 1. k-medoids and average linkage perform similar, as they heavily rely on the computation of pairwise distances and thus are sensitive to noise, and cannot well capture the complex cluster structures in the real data sets. NCuts, NJW-SC, and Zell have good performance on most data sets, as they capture the underlying manifold structures of the data. STSC works fine on some synthetic multiscale datasets in [11] but its results are worse than ours on several real datasets in comparison. Note that STSC adaptively estimated the parameter σ^2 at every point to reflect the variation of local density while ours explores indgree/outdegree and fixes σ^2 as constant. The effective and robust affinity measure for agglomerative clustering makes our GDL-U and AGDL algorithm performs the best among all the algorithms. The AGDL's results are nearly the same as GDL-U.

Compared to other graph-based algorithms, GDL-U and AGDL are more robust to the parameter σ for building the graph, as well as the noise in the data (see Fig. 4). The noise added to images can degrade the performance of other algorithms greatly, but our performance is barely affected.

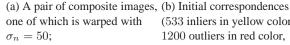
For the graph-based algorithms, we show their time cost in Table 2. AGDL costs the least amount of time among all the algorithms. GDL is faster than NCuts, NJW-SC, and DGSC, and is much faster than Zell. G-Link, which has worse performance than AGDL, is comparable to AGDL on time cost.

Dataset Sample Num Cluster Num NCuts NJW-SC DGSC Zell GDL-U AGDL 15.22 COIL-20 1440 20 3.880 6.399 8.280 0.265 0.277 COIL-100 7200 100 133.8 239.7 326.4 432.9 12.81 5.530 USPS 11000 10 263.0 461.6 538.9 9703 53.64 29.01 247.2 384 4 35 60 17.18 MNIST 10000 10 460 4 64003 Yale-B 2414 38 9.412 13.99 16.00 178.2 0.731 0.564 FRGC 222 12776 577 4 9143 1012.2 65021 49 15 18.62

Table 2. The time cost (in *seconds*) of the algorithms. The minimum time cost on each dataset is in hold. The statistics of each dataset are shown for reference.

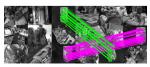








(533 inliers in vellow color. 1200 outliers in red color. according to the ground truth):



(c) Detected correspondences by AGDL (532 true, 552 detected, F-score 0.981).

Fig. 5. Example of object matching through feature correspondence clustering.

4.2 **Feature Correspondence Clustering for Object Matching**

We show the effectiveness of our clustering algorithm in the presence of outliers via feature correspondence clustering. Feature correspondence clustering is commonly used for robust object matching [21,14,22], which can deal with geometric distortions of objects across images and is a fundamental problem in computer vision. We demonstrate that our algorithm can be effectively integrated with the framework of feature correspondence clustering. Therefore, it has a range of potential applications, such as object recognition, image retrieval, and 3D reconstruction.

We compare with two recent state-of-the-art methods, i.e., agglomerative correspondence clustering (ACC) [14] and graph shift (GS) [22].⁵

Overview of experiments. We follow the experiments in the ACC paper [14]. We use composite images and their warped versions (Fig. 5(a)) to simulate cluttered scenes where deformable objects appear. Then we can use the ground-truth for performance evaluation. Namely, we compute the precision and recall rates of detected correspondences (Fig. 5(c)), given a set of correspondences with ground-truth (Fig. 5(b)). A good clustering algorithm can group inliers and separate outliers. It is a more direct way of evaluating the performance of clustering algorithms than other experiments, such as object recognition.

⁵ The code of ACC and GS are downloaded from http://cv.snu.ac.kr/research/~acc/, and http://sites.google.com/site/lhrbss/, respectively. We do not present the results of spectral matching (SM) [21], because both ACC and GS outperformed SM greatly [14,22], especially when there existed at least two clusters of correspondences according to the ground-truth.

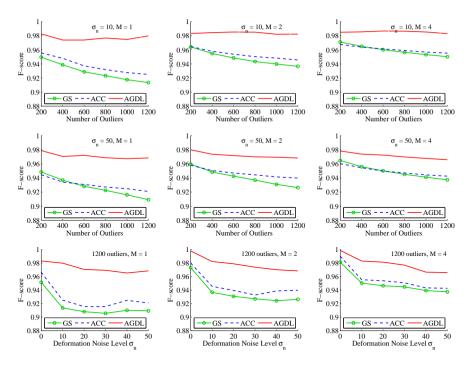


Fig. 6. Performance comparison of different algorithms. In each sub-figure, one of the three factors, i.e., the number of outliers, the level of deformation σ_n , and the number of common sub-images M, is varied, while the other two are fixed as the values appearing at the top. All the results are averaged over 30 random trials.

Experimental settings. We generate a pair of 3×3 tiled images that contain M common sub-image(s). The common sub-images are randomly selected from the model images of the ETHZ toys dataset⁶, and the non-common sub-images are from test images of the same dataset. The positions of all sub-images are randomly determined. When M > 1, the common sub-images are chosen as different objects. To simulate deformation, one of the paired images is warped using the thin-plate spline (TPS) model. An example of paired test images are shown in Fig. 5(a), 9×9 crossing points from a 10×10 meshgrid on the image are chosen as the control points of the TPS model. Then, all the control points are perturbed by Gaussian noise of $N(0, \sigma_n^2)$ independently, and the TPS warping is applied based on the perturbations of control points. To obtain the candidate correspondences between two tiled images, features are extracted by the MSER detector, and the best 3,000 correspondences are collected according to similarity of the SIFT descriptors. Using the warping model, each correspondence has a ground-truth label: true if its error is smaller than three pixels, and false otherwise. Fig. 5(b) shows the correspondences as lines, among which the yellow ones represent true correspondences. Then, the performance of different algorithms are quantitatively eval-

 $^{^6}$ http://www.vision.ee.ethz.ch/~calvin/datasets.html.

uated. We use the F-score, a traditional statistical measure of accuracy, which is defined as $[precision \cdot recall/(precision + recall)]$.

Parameters of ACC and GS. As we strictly follow the test protocol in the ACC paper [14], we use the default parameters in their codes. For GS, we compute the affinity matrix $\mathbf{W}_{ij} = \max(\beta - d_{ij}/\sigma_s^2, 0)$ as the paper [22], where d_{ij} is the distance between correspondence i and correspondence j as defined in the ACC paper [14]. β , σ_s and other parameters in GS are tuned to be the best.

Parameters of AGDL. For our AGDL algorithm (i.e., Algorithm 3), the parameters are fixed as $n_T = 50$, a = 10, K = 35, $K^0 = 2$, and $K^c = 10$. We found that the GDL works well in a large range of n_T , as the number of ground-truth clusters (i.e., M) is very small and we can eliminate the outlier clusters by postprocessing.⁷

Results. As shown in Fig. 6, we vary the number of outliers, the level of deformation, and the number of common sub-images, and compare the F-scores of detected correspondences by different algorithms. Both ACC and GS perform excellently on this task. It is challenging to beat them, which are very recent methods designed specifically for object matching. However, our simple clustering algorithm outperforms them. We find our AGDL algorithm performs consistently better than both ACC and GS under different settings. AGDL has a higher F-score than both in 95.6% of the random trials under all the setting combinations. We attribute the success of AGDL to the effective cluster affinity measure which is robust to noise and outliers.

5 Conclusion

We present a fast and effective method for agglomerative clustering on a directed graph. Our algorithm is based on indegree and outdegree, fundamental concepts in graph theory. The indegree and outdegree have been widely studied in complex networks, but have not received much attention in clustering. We analyze their roles in modeling the structures of data, and show their power via the proposed graph degree linkage algorithm. We demonstrated the superiority of this simple algorithm on image clustering and object matching. We believe our work provides not only a simple and powerful clustering algorithm to many applications in computer vision, but also an insightful analysis of the graph representation of data via indegree and outdegree.

Acknowledgment

This work is partially supported by the General Research Fund sponsored by the Research Grants Council of Hong Kong (Project No. CUHK416510, CUHK417110 and CUHK417011) and National Natural Science Foundation of China (Project No.61005057). It is also supported through Introduced Innovative R&D Team of Guangdong Province 201001D0104648280 and Shenzhen Key Lab of Computer Vision and Pattern Recognition. The authors would like to thank Tianfan Xue for proof reading and Wei Li for help on the ACC code.

⁷ Please see Sec. 9 in the supplemental materials for details of outlier elimination in our algorithm. Different from ACC, which utilizes *additional* information, i.e., geometrical locations of feature points, we only use the *K*-NN graph in outlier elimination.

References

- 1. Hastie, T., Tibshirani, R., Friedman, J.: The elements of statistical learning: Data mining, inference, and prediction, second edn. Springer Verlag (2009)
- Shi, J., Malik, J.: Normalized cuts and image segmentation. IEEE TPAMI 22(8) (2000) 888–905
- 3. Ng, A., Jordan, M., Weiss, Y.: On spectral clustering: Analysis and an algorithm. In: NIPS. (2001)
- Belkin, M., Niyogi, P.: Laplacian eigenmaps for dimensionality reduction and data representation. Neural Computation 15(6) (2003) 1373–1396
- Grady, L., Schwartz, E.: Isoperimetric graph partitioning for image segmentation. IEEE TPAMI 28(3) (2006) 469–475
- Zhang, W., Lin, Z., Tang, X.: Learning semi-Riemannian metrics for semisupervised feature extraction. IEEE TKDE 23(4) (2011) 600–611
- 7. Frey, B., Dueck, D.: Clustering by passing messages between data points. Science **315**(5814) (2007) 972–976
- 8. Zhou, D., Huang, J., Schölkopf, B.: Learning from labeled and unlabeled data on a directed graph. In: ICML. (2005)
- Kleinberg, J.: Authoritative sources in a hyperlinked environment. Journal of the ACM 46(5) (1999) 604–632
- Mislove, A., Marcon, M., Gummadi, K., Druschel, P., Bhattacharjee, B.: Measurement and analysis of online social networks. In: Proc. ACM SIGCOMM Conf. on Internet Measurement. (2007)
- 11. Zelnik-Manor, L., Perona, P.: Self-tuning spectral clustering. In: NIPS. (2005)
- 12. Wu, M., Schölkopf, B.: A local learning approach for clustering. In: NIPS. (2007)
- 13. Franti, P., Virmajoki, O., Hautamaki, V.: Fast agglomerative clustering using a k-nearest neighbor graph. IEEE TPAMI **28**(11) (2006) 1875–1881
- Cho, M., Lee, J., Lee, K.: Feature correspondence and deformable object matching via agglomerative correspondence clustering. In: ICCV. (2009)
- Sander, J., Ester, M., Kriegel, H., Xu, X.: Density-based clustering in spatial databases: The algorithm GDBSCAN and its applications. Data Mining and Knowledge Discovery 2(2) (1998) 169–194
- Ertöz, L., Steinbach, M., Kumar, V.: Finding clusters of different sizes, shapes, and densities in noisy, high dimensional data. In: SIAM International Conf. on data mining. (2003)
- Karypis, G., Han, E., Kumar, V.: Chameleon: Hierarchical clustering using dynamic modeling. IEEE Computer 32(8) (1999) 68–75
- 18. Zhao, D., Tang, X.: Cyclizing clusters via zeta function of a graph. In: NIPS. (2008)
- Felzenszwalb, P., Huttenlocher, D.: Efficient graph-based image segmentation. IJCV 59(2) (2004) 167–181
- 20. Yu, S., Shi, J.: Multiclass spectral clustering. In: ICCV. (2003)
- Leordeanu, M., Hebert, M.: A spectral technique for correspondence problems using pairwise constraints. In: ICCV. (2005)
- Liu, H., Yan, S.: Common visual pattern discovery via spatially coherent correspondences. In: CVPR. (2010)

6 Implementations of GDL

Algorithm 2 Graph Degree Linkage with the update formula (GDL-U)

Input: a set of n samples $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n\}$, and the target number of clusters n_T . Build the K^0 -NN graph, and detect its weakly connected components as initial clusters. Denote the set of initial clusters as $V^c = \{\mathcal{C}_1, \cdots, \mathcal{C}_{n_c}\}$, where n_c is the number of clusters.

Build the K-NN graph, and get the weighted adjacency matrix \mathbf{W} .

Initialize the asymmetric affinity table $\mathcal{A}_{\mathcal{C}_a \to \mathcal{C}_b}$ for $\mathcal{C}_a, \mathcal{C}_b \in V^c$.

while $n_c > n_T$ do

```
Search two clusters C_a and C_b, such that \{C_a, C_b\} = \operatorname{argmax}_{C_a, C_b \in V^c} A_{C_a, C_b};
V^c \leftarrow \{V^c \setminus \{C_a, C_b\}\} \cup \{C_{ab}\}, \text{ where } C_{ab} = C_a \cup C_b, \text{ and } n_c = n_c - 1;
```

For all C_c , compute $A_{C_{ab} \to C_c}$ using the update formula, i.e., Eq. (9), and $A_{C_c \to C_{ab}}$ using Eq. (8).

end while

Output: V^c .

Algorithm 3 Accelerated Graph Degree Linkage (AGDL)

Input: a set of n sample vectors $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n\}$, and the target number of clusters n_T .

Build the K^0 -NN graph, and detect its weakly connected components as initial clusters. Denote the set of initial clusters as $V^c = \{C_1, \dots, C_{n_c}\}$, where n_c is the number of clusters.

Build the K-NN graph, and get the weighted adjacency matrix \mathbf{W} .

Create a neighbor set for each cluster in V^c , and initialize it as the K^c -nearest cluster set.

```
while n_c > n_T do
```

Search two clusters \mathcal{C}_a and \mathcal{C}_b from the affinity of pairs of clusters associated with the neighbor sets, such that $\{\mathcal{C}_a, \mathcal{C}_b\} = \operatorname{argmax}_{\mathcal{C}_a \in \mathcal{N}_{\mathcal{C}_b}}$ or $\mathcal{C}_b \in \mathcal{N}_{\mathcal{C}_a}$ $\mathcal{A}_{\mathcal{C}_a, \mathcal{C}_b}$;

 $V^c \leftarrow \{V^c \setminus \{C_a, C_b\}\} \cup \{C_{ab}\}, \text{ where } C_{ab} = C_a \cup C_b, \text{ and } n_c = n_c - 1;$

For all C_c , such that $C_a \in \mathcal{N}_{C_c}$ or $C_b \in \mathcal{N}_{C_c}$, add C_{ab} to \mathcal{N}_{C_c} , and compute the affinity \mathcal{A}_{C_{ab},C_c} ;

Find the K^c -nearest clusters for \mathcal{C}_{ab} in the set $\mathcal{N}_{\mathcal{C}_a} \cup \mathcal{N}_{\mathcal{C}_b}$, to form $\mathcal{N}_{\mathcal{C}_{ab}}$;

Remove C_a and C_b from the neighbor sets, and remove \mathcal{N}_{C_a} and \mathcal{N}_{C_b} .

end while

Output: V^c .

7 Proof of Theorem 3

Proof. For (a), we analyze the time complexity for each part of the GDL algorithm.

(1) The directed graph construction has a complexity of at most $O(Kn^2)$ (naive implementation). Note that $n \gg K$, and thus we omit K in the complexities hereinafter.

- (2) The complexity of constructing initial clusters is O(n), as the number of edges in the graph G is $O(K^0n)$, where K^0 is 1 or 2.
- (3) In our clustering algorithm, we use an $n_c \times n_c$ table to store the affinities between clusters. As the initial clusters are of small sizes, we can assume O(1) complexity for computing the affinity between each pair of two clusters. So, it requires a complexity of $O(n_0^2)$ to initialize the table, where n_0 is the number of initial clusters $(n_0 < n)$.
- (4) In each iteration, it costs $O(n_c^2)$ to find the maximum value in the cluster affinity table. To update the cluster affinity table after merging the two clusters with maximum affinity value, we need to compute $(n_c 1)$ affinities, and each affinity is computed with complexity of $O(|\mathcal{C}_a| + |\mathcal{C}_b|)$ using Eq. (5) (because W is a sparse matrix with K nonzero elements in each row). Therefore, the complexity for each iteration is at most $O(n_c n)$.
- (5) The number of iterations is $(n_0 n_T)$.

By replacing n_0 and n_c with their upper bound n, a loose upper bound of the time complexity of the GDL algorithm is $O(n^3)$.

For (b), we can reduce the complexity in each iteration from $O(n_c n)$ in (a) to O(n). We can maintain a table to store the nearest cluster of each cluster.⁸ In each iteration, finding the maximum value and updating the table cost approximately $O(n_c)$. For the affinity table, the updating scheme of $\mathcal{A}_{\mathcal{C}_{ab} \to \mathcal{C}_c}$ as in Eq. 9 costs $O(n_c)$ for all the new affinities. To compute $\mathcal{A}_{\mathcal{C}_c \to \mathcal{C}_{ab}}$, the total complexity for all the new affinities is less than the complexity of computing $\mathbf{W}_{\mathcal{C}_{ab},*}\mathbf{W}_{*,\mathcal{C}_{ab}}$, which is O(nK), as $\mathbf{W}_{\mathcal{C}_{ab},*}$ is K-sparse in each row. $\mathbf{W}_{\mathcal{C}_{ab},*}$ is the submatrix of \mathbf{W} whose row indices correspond to the vertices in \mathcal{C}_{ab} and column indices are from 1 to n.

Finally, the total complexity for GDL-U is $O(n^2)$.

For (c), there are several differences in the AGDL:

- In (3), we use the neighbor sets of clusters instead of the cluster affinity table. The construction of all the neighbor sets costs $O(K^c n_0^2)$.
- In (4), we need to find the maximum affinity value in the neighbor sets (with complexity of $O(K^c n_c)$ and compute $O(K^c (1+\tau))$ affinities to update the neighbor sets with complexity of $O(K^c n)$). Because the size of the union of neighbor sets of \mathcal{C}_a and \mathcal{C}_b is less than $2K^c$), and for real data, we can assume that the number of clusters whose neighbor set includes \mathcal{C}_a or \mathcal{C}_b is less than $2\tau K^c$, where τ is usually a small constant close to 1. Therefore, the complexity for each iteration is at most O(n).

So, the time complexity of the AGDL algorithm is $O(n^2)$.

⁸ We can use a heap to achieve better efficiency for this part. But it is not the bottleneck for both the complexity analysis and run-time of GDL-U.

Table 3. Quantitative clustering results in CE on real imagery data. A smaller CE value indicates a better clustering result. The results shown in a boldface are significantly better than the others, with a significance level of 0.01.

Dataset	k-med	Link	G-Link	NCuts	NJW-SC	DGSC	STSC	Zell	GDL-U	AGDL
COIL-20	0.401	0.677	0.213	0.246	0.228	0.201	0.158	0.187	0.142	0.142
COIL-100	0.570	0.819	0.394	0.462	0.411	0.396	0.391	0.351	0.267	0.269
USPS	0.607	0.874	0.252	0.459	0.354	0.255	0.421	0.332	0.246	0.246
MNIST	0.577	0.776	0.162	0.405	0.432	0.230	0.305	0.400	0.150	0.150
Yale-B	0.728	0.847	0.376	0.273	0.270	0.237	0.205	0.464	0.197	0.197
FRGC	0.728	0.753	0.664	0.565	0.596	0.595	0.580	0.560	0.548	0.551

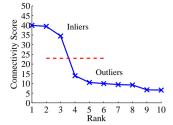


Fig. 7. The connectivity scores of clusters sorted in descending order. The threshold for separating inliers and outliers is shown in a red dash line.

8 Quantitative Results in Clustering Error for Image Clustering

The quantitative results, measured in CE [12], are given in Table 3. The CE is defined as the minimum overall error rate among all possible permutation mappings between true class labels and clusters. A smaller CE value indicates a better clustering result.

9 Outlier Elimination for Object Matching

For AGDL, we observe that there are many inedges and outedges inside a cluster of inliers, while less edges inside a cluster of outliers because outliers are in low density regions. Inspired by this, we define the connectivity score of a cluster \mathcal{C} as

 $\sum_{i \in \mathcal{C}} \left[\deg_i^-(\mathcal{C}) + \deg_i^+(\mathcal{C}) \right]$. We find that there are always large differences between the scores of inlier clusters and outlier clusters (see Fig. 7). Therefore, we rank the final clusters by their connectivity scores. Namely, we sort their scores in descending order, and then search the largest gap between two consecutive scores. The set of clusters are divided into two subsets without intersection. The subset of clusters with small scores is treated as the collection of outliers and removed. For ACC and GS, we use their default methods for outlier elimination.