

1 A new Interactive Semi-Supervised Clustering model for
2 large image database indexing

3 LAI Hien Phuong^{a,b}, Muriel VISANI^a, Alain BOUCHER^{a,b}, Jean-Marc
4 OGIER^a

5 ^a*L3I, Université de La Rochelle, Avenue M. Crépeau, 17042 La Rochelle cedex 1, France*

6 ^b*IFI, Equipe MSI; IRD, UMI 209 UMMISCO*

7 *Institut de la Francophonie pour l'Informatique, 42 Ta Quang Buu, Hanoi, Vietnam*

8 *Vietnam National University, Hanoi, Vietnam*

Abstract

Indexing methods play a very important role in finding information in large image databases. They organize indexed images in order to facilitate, accelerate and improve the results for later retrieval. Alternatively, clustering may be used for structuring the feature space so as to organize the dataset into groups of similar objects without prior knowledge (unsupervised clustering) or with a limited amount of prior knowledge (semi-supervised clustering).

In this paper, we introduce a new interactive semi-supervised clustering model where prior information is integrated via pairwise constraints between images. The proposed method allows users to provide feedback in order to improve the clustering results according to their wishes. Different strategies for deducing pairwise constraints from user feedback were investigated. Our ex-

Email addresses: Corresponding Author: L3I, Université de La Rochelle, Avenue M. Crépeau, 17042 La Rochelle cedex 1, France; Phone: +33 6 46 51 12 32; Fax: +33 5 46 45 82 42; Email: hien_phuong.lai@univ-lr.fr (LAI Hien Phuong), muriel.visani@univ-lr.fr (Muriel VISANI), alainboucher12@gmail.com (Alain BOUCHER), jean-marc.ogier@univ-lr.fr (Jean-Marc OGIER)

periments on different image databases (Wang, PascalVoc2006, Caltech101) show that the proposed method outperforms semi-supervised HMRF-kmeans (Basu et al., 2004).

1 *Keywords:* Semi-supervised clustering, Interactive learning, Image indexing

2 **1. Introduction**

3 Content-Based Image Retrieval (CBIR) refers to the process which uses
4 visual information (usually encoded using color, shape, texture feature vec-
5 tors, etc.) to search for images in the database that correspond to the user's
6 queries. Traditional CBIR systems generally rely on two phases. The first
7 phase is to extract the feature vectors from all the images in the database and
8 to organize them into an efficient index data structure. The second phase
9 is to efficiently search in the indexed feature space to find the most similar
10 images to the query image.

11 With the development of many large image databases, an exhaustive
12 search is generally intractable. Feature space structuring methods (normally
13 called indexing methods) are therefore necessary for facilitating and acceler-
14 ating further retrieval. They can be classified into space partitioning methods
15 and data partitioning methods.

16 Space partitioning methods (KD-tree (Bentley, 1975), KDB-tree (Robin-
17 son, 1981), LSD-tree (Henrich et al., 1989), Grid-File (Nievergelt et al.,
18 1988)...) generally divide the feature space into cells (sometimes referred
19 to as "buckets") of fairly similar cardinality (in terms of number of images
20 per cell), without taking into account the distribution of the images in the
21 feature space. Therefore, dissimilar points may be included in a same cell

1 while similar points may end up in different cells. The resulting index is there-
2 fore not optimal for retrieval, as the user generally wants to retrieve similar
3 images to the query image. Moreover, these methods are not designed to
4 handle high dimensional data, while image feature vectors commonly count
5 hundreds of elements.

6 Data partitioning methods (B-tree (Bayer and McCreight, 1972), R-trees
7 (Guttman, 1984; Sellis et al., 1987; Beckmann et al., 1990), SS-tree (White
8 and Jain, 1996), SR-tree (Katayama and Satoh, 1997), X-tree (Berchtold
9 et al., 1996)...) also integrate information about image distribution in the
10 feature space. However, the limitations on the cardinality of the space cells
11 remain, causing the resulting index to be non-optimal for retrieval, especially
12 in the case where groups of similar objects are unbalanced, *i.e.* composed of
13 different numbers of images.

14 Our claim is that using clustering instead of traditional indexing to or-
15 ganize feature vectors, results in indexes better adapted to high dimensional
16 and unbalanced data. Indeed, clustering aims to split a collection of data
17 into groups (clusters) so that similar objects belong to the same group and
18 dissimilar objects are in different groups, with no constraints on the cluster
19 size. This makes the resulting index better optimized for retrieval. In fact,
20 while in traditional indexing methods it might be difficult to fix the number
21 of objects in each bucket (especially in the case of unbalanced data), cluster-
22 ing methods have no limitation on the cardinality of the clusters, objects can
23 be grouped into clusters of very different sizes. Moreover, using clustering
24 might simplify the relevance feedback task, as the user might interact with
25 a small number of cluster prototypes rather than numerous single images.

1 Because feature vectors only capture low level information such as color,
2 shape or texture, there is a semantic gap between high-level semantic con-
3 cepts expressed by the user and these low-level features. The clustering
4 results are therefore generally different from the intent of the user. Our work
5 aims to involve users in the clustering phase so that they can interact with
6 the system in order to improve the clustering results. The clustering meth-
7 ods should therefore produce a hierarchical cluster structure where the initial
8 clusters may be easily merged or split. We are also interested in clustering
9 methods which can be incrementally built in order to facilitate the insertion
10 or deletion of new images by the user. It can be noted that incrementality is
11 also very important in the context of huge image databases, when the whole
12 dataset cannot be stored in the main memory. Another very important point
13 is the computational complexity of the clustering algorithm, especially in an
14 interactive online context where the user is involved.

15 In the case of large image database indexing, we may be interested in tra-
16 ditional clustering (unsupervised) (Jain et al., 1999; Xu and Wunsch, 2005)
17 or semi-supervised clustering (Basu et al., 2002; Dubey et al., 2010; Wagstaff
18 et al., 2001; Basu et al., 2004). While no information about ground truth is
19 provided in the case of unsupervised clustering, a limited amount of knowl-
20 edge is available in the case of semi-supervised clustering. The provided
21 knowledge may consist of class labels (for some objects) or pairwise con-
22 straints (must-link or cannot-link) between objects.

23 In (Lai et al., 2012a), we proposed a survey of unsupervised clustering
24 techniques and analyzed the advantages and disadvantages of different meth-
25 ods in a context of huge masses of data where incrementality and hierarchi-

1 cal structuring are needed. We also experimentally compared five methods
2 (global k-means (Likas et al., 2003), AHC (Lance and Williams, 1967), R-tree
3 (Guttman, 1984), SR-tree (Katayama and Satoh, 1997) and BIRCH (Zhang
4 et al., 1996)) with different real image databases of increasing sizes (Wang,
5 PascalVoc2006, Caltech101, Corel30k) (the number of images ranges from
6 1,000 to 30,000) to study the scalability of different approaches relative to
7 the size of the database. In (Lai et al., 2012b), we presented an overview of
8 semi-supervised clustering methods and proposed a preliminary experiment
9 of an interactive semi-supervised clustering model using the HMRF-kmeans
10 (Hidden Markov Random Fields kmeans) clustering (Basu et al., 2004) on the
11 Wang image database in order to analyze the improvement in the clustering
12 process when user feedback is provided.

13 There are three main parts to this paper. Firstly, we propose a new inter-
14 active semi-supervised clustering model using pairwise constraints. Secondly,
15 we investigate different methods for deducing pairwise constraints from user
16 feedback. Thirdly, we experimentally compare our proposed semi-supervised
17 method with the widely known semi-supervised HMRF-kmeans method.

18 This paper is structured as follows. A short review of semi-supervised
19 clustering methods is presented in Section 2. Our interactive semi-supervised
20 clustering model is proposed in Section 3. Some experiments are presented
21 in Section 4. Some conclusions and further works are provided in Section 5.

22 **2. A short review of semi-supervised clustering methods**

23 For unsupervised clustering only similarity information is used to orga-
24 nize objects; in the case of semi-supervised clustering a small amount of prior

1 knowledge is available. Prior knowledge is either in the form of class labels
2 (for some objects) or pairwise constraints between objects. Pairwise con-
3 straints specify whether two objects should be in the same cluster (must-link)
4 or in different clusters (cannot-link). As the clusters produced by unsuper-
5 vised clustering may not be the ones required by the user, this prior knowl-
6 edge is needed to guide the clustering process for resulting clusters which are
7 closer to the user’s wishes. For instance, for clustering a database with thou-
8 sands of animal images, an user may want to cluster by animal species or by
9 background landscape types. An unsupervised clustering method may give,
10 as a result, a cluster containing images of elephants with a grass background
11 together with images of horses with a grass background and another cluster
12 containing images of elephants with a sand background. These results are
13 ideal when the user wants to cluster by background landscape types. But
14 they are poor when the user wants to cluster by animal species. In this case,
15 must-link constraints between images of elephants with a grass background
16 and images of elephants with a sand background and cannot-link constraints
17 between images of elephants with a grass background and images of horses
18 with a grass background are needed to guide the clustering process. The
19 objective of our work is to make the user interact with the system so as to
20 define easily these constraints with only a few clicks. Note that the avail-
21 able knowledge is too poor to be used with supervised learning, as only a
22 very limited ratio of the available images are considered by the user at each
23 step. In general, semi-supervised clustering methods are used to maximize
24 intra-cluster similarity, to minimize inter-cluster similarity and to keep a high
25 consistency between partitioning and domain knowledge.

1 Semi-supervised clustering has been developed in the last decade and
2 some methods have been published to date. They can be divided into semi-
3 supervised clustering with labels, where partial information about object
4 labels is given, and semi-supervised clustering with constraints, where a small
5 amount of pairwise constraints between objects is given.

6 Some semi-supervised clustering methods using labeled objects have been
7 put forward: seeded-kmeans (Basu et al., 2002), constrained-kmeans (Basu
8 et al., 2002), etc. Seeded-kmeans and constrained-kmeans are based on the k-
9 means algorithm. Prior knowledge for these two methods is a small subset of
10 the input database, called seed set, containing user-specified labeled objects
11 of k different clusters. Unlike k-means algorithm which randomly selects
12 the initial cluster prototypes, these two methods use the labeled objects to
13 initialize the cluster prototypes. Following this we repeat, until convergence,
14 the re-assignment of each object in the dataset to the nearest prototype
15 and the re-computation of the prototypes with the assigned objects. The
16 seeded-kmeans assigns objects to the nearest prototype without considering
17 the prior labels of the objects in the seed set. In contrast, the constrained-
18 kmeans maintains the labeled examples in their initial clusters and assigns
19 the other objects to the nearest prototype. An interactive cluster-level semi-
20 supervised clustering was proposed in (Dubey et al., 2010) for document
21 analysis. In this model, knowledge is progressively provided as assignment
22 feedback and cluster description feedback after each interactive iteration.
23 Using assignment feedback, the user moves an object from one cluster to
24 another cluster. Using cluster description feedback, the user modifies the
25 feature vector of any current cluster (*e.g.* increase the weighting of some

1 important words). The algorithm learns from all the feedback to re-cluster
 2 the dataset in order to minimize average distance between points and their
 3 cluster centers while minimizing the violation of constraints corresponding
 4 to feedback.

5 Among the semi-supervised clustering methods using pairwise constraints
 6 between objects, we can cite COP-kmeans (constrained-kmeans) (Wagstaff
 7 et al., 2001), HMRF-kmeans (Hidden Markov Random Fields Kmeans) (Basu
 8 et al., 2004), semi-supervised kernel-kmeans (Kulis et al., 2005), etc. The
 9 input data of these methods is data set X , a set of must-link constraints M
 10 and a set of cannot-link constraints C . In COP-kmeans, points are assigned
 11 to clusters without violating any constraint. A point x_i is assigned to its
 12 closest cluster μ_j unless a constraint is violated. If x_i cannot be placed in
 13 μ_j , we continue attempting to assign x_i to the next cluster in the sorted list
 14 of clusters by ascending order of distances with x_i until a suitable cluster
 15 is found. The clustering fails if no solution respecting the constraints is
 16 found. While the constraint violation is strictly prohibited in COP-kmeans,
 17 it is allowed with a violation cost (penalty) in HMRF-kmeans and in semi-
 18 supervised kernel-kmeans. The objective function to be minimized in the
 19 semi-supervised HMRF-kmeans is as follows:

$$J_{HMRF-Kmeans} = \sum_{x_i \in X} D(x_i, \mu_{l_i}) + \sum_{(x_i, x_j) \in M, l_i \neq l_j} w_{ij} + \sum_{(x_i, x_j) \in C, l_i = l_j} \overline{w_{ij}} \quad (1)$$

20 where w_{ij} ($\overline{w_{ij}}$) is the penalty cost for violating a must-link (cannot-link)
 21 constraint between x_i and x_j , l_i refers to the cluster label of x_i , and $D(x_i, \mu_{l_i})$
 22 measures the distance between x_i and its corresponding cluster center μ_{l_i} .
 23 The violation cost of a pairwise constraint may be either a constant or a

1 function of the distance between the two points specified in the pairwise
2 constraint as follows:

$$w_{ij} = wD(x_i, x_j) \quad (2)$$

3

$$\overline{w}_{ij} = \overline{w}(D_{max} - D(x_i, x_j)) \quad (3)$$

4 where w and \overline{w} are constants specifying the cost for violating a must-link or
5 a cannot-link constraint. D_{max} is the maximal distance between two points
6 in the data set. We can see that, to ensure the most difficult constraints are
7 respected, higher penalties are assigned to violations of must-link constraints
8 between points which are distant and to violations of cannot-link constraints
9 between points which are close. The term D_{max} in Equation (3) can make
10 the cannot-link penalty term sensitive to extreme outliers, but all cannot-link
11 constraints are treated in the same way, so even in the presence of extreme
12 outliers, there would be no cannot-link constraint favored compared to the
13 others. The objective function in Equation (1) is also sensitive to outliers.
14 We can reduce this sensitivity by using an outlier filtering technique or by
15 replacing the term D_{max} by the maximum distance between two clusters.
16 HMRF-kmeans first initializes the k cluster centers based on user-specified
17 constraints, as described in (Basu et al., 2004). After the initialization step,
18 an iterative relocation approach similar to k-means is applied to minimize the
19 objective function. The iterative algorithm represents the repetition of the
20 assignment phase of each point to the cluster which minimizes its contribution
21 to the objective function and the re-estimation phase of the cluster centers
22 minimizing the objective function. The semi-supervised kernel-kmeans (Kulis
23 et al., 2005) is similar to the HMRF-kmeans, but calculates the objective
24 function in a transformed space instead of the original space using a kernel

1 function mapping as follows:

$$J_{SS_Kernel_Kmeans} = \sum_{x_i \in X} \|\phi(x_i) - \bar{\phi}_{l_i}\|^2 - \sum_{(x_i, x_j) \in M, l_i = l_j} w_{ij} + \sum_{(x_i, x_j) \in C, l_i \neq l_j} \bar{w}_{ij} \quad (4)$$

2 where $\phi(x_i)$ is the kernel function mapping, $\bar{\phi}_{l_i}$ is the centroid of the cluster containing x_i and w_{ij} (\bar{w}_{ij}) is the penalty cost for violating a must-link (cannot-link) constraint between x_i and x_j . In the second term of Equation 3 (4), instead of adding a penalty cost for a must-link violation if the two points 4 are in different clusters, Kulis et al. (2005) give a reward for must-link constraint satisfaction if the two points are in the same cluster, by subtracting 5 the corresponding penalty term from the objective function. 6 7 8

9 **3. Proposed interactive semi-supervised clustering model**

10 In this section, we present our proposed interactive semi-supervised clustering model. In our model, the initial clustering is carried out without any 11 prior knowledge, using an unsupervised clustering method. In Lai et al. 12 (2012a) we discussed the adequation between different unsupervised clustering methods and our applied context (involving user interactivity) as 13 well as experimentally compared different unsupervised clustering methods 14 (global k-means (Likas et al., 2003), AHC (Lance and Williams, 1967), R-tree (Guttman, 1984), SR-tree (Katayama and Satoh, 1997), BIRCH (Zhang 15 et al., 1996)). Our conclusion was that BIRCH is the most suitable to our 16 context. BIRCH is less sensitive to variations in its parameters. Moreover, 17 it is incremental, it provides a hierarchical structure of clusters and it out- 18 performs other methods in the context of a large database (best results and 19 best computational time in our tests). Therefore, BIRCH is chosen for the 20 21 22

1 initial unsupervised clustering in our model. After the initial clustering, the
 2 user views the clustering results and provides feedback to the system. The
 3 pairwise constraints (must-link, cannot-link) are deduced, based on user feed-
 4 back; the system then re-organizes the clusters by considering the constraints.
 5 The re-clustering process is done using the proposed semi-supervised cluster-
 6 ing described in Section 3.2. The interactive process (user provides feedback
 7 and system reorganizes the clusters) is repeated until the clustering result
 8 satisfies the user. The interactive semi-supervised clustering model contains
 9 the following steps:

- 10 1. Initial clustering using BIRCH unsupervised clustering.
- 11 2. *Repeat*:
 - 12 (a) Receive feedback from the user and deduce pairwise constraints.
 - 13 (b) Re-organize the clusters using the proposed semi-supervised cluster-
 14 tering method.
- 15 *until* the clustering result satisfies the user.

16 3.1. BIRCH unsupervised clustering

17 Let us briefly describe the BIRCH (Balanced Iterative Reducing and
 18 Clustering using Hierarchies) unsupervised clustering method (Zhang et al.,
 19 1996). The idea of BIRCH is to build a Clustering Feature Tree (CF-tree).

20 We define a CF-vector, summarizing information of a cluster including N
 21 vectors $(\vec{x}_1, \dots, \vec{x}_N)$, as a triplet $CF = (N, \vec{LS}, SS)$, where \vec{LS} and SS are
 22 respectively the linear sum and the square sum of vectors ($\vec{LS} = \sum_{i=1}^N \vec{x}_i$;
 23 $SS = \sum_{i=1}^N \vec{x}_i^2$). From the CF-vectors, we can simply compute the centroid,

1 the radius (average distance from points to the centroid) of a cluster and also
 2 the distance between two clusters (*e.g.* the Euclidean distance between their
 3 centroids). A CF-tree is a balanced tree having three parameters B , L and
 4 T :

- 5 • Each internal node contains, at most, B elements of the form $[CF_i, child_i]$
 6 where $child_i$ is a pointer to its i^{th} child node and CF_i is the CF-vector
 7 of this child.
- 8 • Each leaf node contains, at most, L entries of the form $[CF_i]$, it also
 9 contains two pointers, $prev$ and $next$, to link leaf nodes.
- 10 • Each entry CF_i represents the information of a group of points which
 11 are close together. Each entry CF_i of a leaf node must have a radius
 12 lower than a threshold T (threshold condition).

13 The CF-tree is created by successively inserting points into the tree. A
 14 new point is preferably inserted in the closest CF_i of the closest leaf, if the
 15 threshold condition is not violated. If it is impossible, a new CF_j is created
 16 for the new point. The corresponding internal and leaf nodes must be split
 17 if necessary. After creating the CF-tree, we can use any clustering method
 18 (AHC, k-means, etc.) to cluster all leaf entries CF_i . In our work, we use
 19 k-means for clustering the leaf entries, as it is suitable to be used with our
 20 proposed semi-supervised clustering in the interactive phase.

21 3.2. Proposed semi-supervised clustering method

22 At each interactive iteration, our semi-supervised clustering method is
 23 applied after receiving feedback from the users for re-organizing the clusters

1 according to their wishes. Our semi-supervised clustering method considers
 2 the set of all leaf entries $S_{CF} = (CF_1, \dots, CF_m)$ of the CF-tree. Supervised
 3 information is provided as two sets of pairwise constraints between CF entries
 4 deduced from user feedback: must-links $M_{CF} = \{(CF_i, CF_j)\}$ and cannot-
 5 links $C_{CF} = \{(CF_i, CF_j)\}$. $(CF_i, CF_j) \in M_{CF}$ implies that CF_i , CF_j and
 6 therefore all points which are included in these two entries should belong to
 7 the same cluster, while $(CF_i, CF_j) \in C_{CF}$ implies that CF_i and CF_j should
 8 belong to different clusters. The objective function to be minimized is as
 9 follows:

$$\begin{aligned}
 J_{obj} &= \sum_{CF_i \in S_{CF}} D(CF_i, \mu_{l_i}) \\
 &+ \sum_{(CF_i, CF_j) \in M_{CF}, l_i \neq l_j} w N_{CF_i} N_{CF_j} D(CF_i, CF_j) \\
 &+ \sum_{(CF_i, CF_j) \in C_{CF}, l_i = l_j} \bar{w} N_{CF_i} N_{CF_j} (D_{max} - D(CF_i, CF_j)) \quad (5)
 \end{aligned}$$

10 where:

- 11 • The first term measures the distortion between each leaf entry CF_i and
 12 the corresponding cluster center μ_{l_i} , l_i refers to the cluster label of CF_i .
- 13 • The second and the third terms represent the penalty costs for re-
 14 spectively violating the must-link and cannot-link constraints between
 15 CF entries. w and \bar{w} are constants specifying the violation cost of
 16 a must-link and a cannot-link between two points. As an entry CF_i
 17 represents the information of a group of N_{CF_i} points, a pairwise con-
 18 straint between two entries CF_i and CF_j corresponds to $N_{CF_i} \times N_{CF_j}$
 19 constraints between points of these two entries. The violation cost of

1 a pairwise constraint between two entries CF_i, CF_j is thus a function
 2 of their distance $D(CF_i, CF_j)$ and of the number of points included
 3 in these two entries. D_{max} is the maximum distance between two CF
 4 entries in the data set. Therefore, higher penalties are assigned to vio-
 5 lations of must-link between entries that are distant and of cannot-link
 6 between entries which are close. As in HMRf-kmeans, the term D_{max}
 7 can make the cannot-link penalty term sensitive to extreme outliers,
 8 and could be replaced by the maximum distance between two clusters
 9 if the database contains extreme outliers.

10 In our case, we use the most frequently used squared Euclidean distance as
 11 distortion measure. The distance between two entries $CF_i = (N_{CF_i}, \vec{LS}_{CF_i},$
 12 $SS_{CF_i}), CF_j = (N_{CF_j}, \vec{LS}_{CF_j}, SS_{CF_j})$ is calculated as the distance between
 13 their means as follows:

$$D(CF_i, CF_j) = \sum_{p=1}^d \left(\frac{LS_{CF_i}(p)}{N_{CF_i}} - \frac{LS_{CF_j}(p)}{N_{CF_j}} \right)^2 \quad (6)$$

14 where d is the number of dimensions of the feature space.

15 The proposed semi-supervised clustering is as follows:

16 **Input:** Set of leaf entries $S_{CF} = \{CF_i\}_{i=1}^m$ which are clustered into K clusters
 17 with the corresponding centroids $\{\mu_h\}_{h=1}^K$,

18 set of must-link constraints $M_{CF} = \{(CF_i, CF_j)\}$

19 set of cannot-link constraints $C_{CF} = \{(CF_i, CF_j)\}$.

20 **Output:** New disjoint K clusters of S_{CF} such that the objective function in
 21 Equation (5) is locally minimized.

22 **Method:**

23 1. Set $t \leftarrow 0$

- 1 2. Repeat until convergence
- 2 (a) Re-assignment step: Given $\{\mu_h^{(t)}\}_{h=1}^K$, re-assign cluster labels $\{l_i^{(t+1)}\}_{i=1}^m$
- 3 of entries $\{CF_i\}_{i=1}^m$ to minimize the objective function.
- 4 (b) Re-estimation step: Given cluster labels $\{l_i^{(t+1)}\}_{i=1}^m$, re-calculate
- 5 the cluster centroids $\{\mu_h^{(t+1)}\}_{h=1}^K$ to minimize the objective func-
- 6 tion.
- 7 (c) $t \leftarrow t + 1$.

8 In the re-assignment step, given the current cluster centers, each entry

9 CF_i is re-assigned to the cluster μ_h which minimizes its contribution to the

10 objective function as follows:

$$\begin{aligned}
J_{obj}(CF_i, \mu_h) &= D(CF_i, \mu_h) \\
&+ \sum_{(CF_i, CF_j) \in M_{CF}, h \neq l_j} w N_{CF_i} N_{CF_j} D(CF_i, CF_j) \\
&+ \sum_{(CF_i, CF_j) \in C_{CF}, h = l_j} \bar{w} N_{CF_i} N_{CF_j} (D_{max} - D(CF_i, CF_j)) \quad (7)
\end{aligned}$$

11 We can see that the optimal assignment of each CF entry also depends on

12 the current assignment of the other CF entries due to the violation cost of

13 pairwise constraints in the second and third terms of Equation 7. Therefore,

14 after all entries are re-assigned, they are randomly re-ordered, and the re-

15 assignment process is repeated until no CF entry changes its cluster label

16 between two successive iterations.

17 In the re-estimation step, given the cluster labels $\{l_i^{(t+1)}\}_{i=1}^m$ of all CF

18 entries, the cluster centers $\{\mu_h\}_{h=1}^K$ are re-calculated in order to minimize

19 the objective function of the current assignment. For simple calculation,

1 each cluster center is also represented in the form of a CF-vector. By using
 2 the squared Euclidean measure, the CF-vector of each cluster prototype μ_h is
 3 calculated based on CF entries which are assigned to this cluster as follows:

$$N_{\mu_h} = \sum_{l_i=h} N_{CF_i} \quad (8)$$

$$\vec{LS}_{\mu_h} = \sum_{l_i=h} \vec{LS}_{CF_i} \quad (9)$$

$$SS_{\mu_h} = \sum_{l_i=h} SS_{CF_i} \quad (10)$$

4 We can see that in each re-assignment step, each entry CF_i moves to
 5 a new cluster μ_h if its contribution to the objective function is decreased
 6 with this re-assignment. Therefore, the objective function J_{obj} is decreased
 7 or unchanged after the re-assignment step. And in each re-estimation step,
 8 the mean of the CF-vector of each cluster μ_h corresponds to the mean of
 9 the CF entries (and therefore the points) in this cluster, that minimizes
 10 the contribution of μ_h to the component $\sum_{CF_i \in S_{CF}} D(CF_i, \mu_{l_i})$ of J_{obj} . The
 11 penalty terms of J_{obj} are not functions of the centroid, thus they do not
 12 take part in cluster center re-estimation. Therefore, the objective function
 13 J_{obj} will decrease or remain the same in the re-estimation step. Since J_{obj}
 14 is bounded below and decreases after each re-assignment and re-estimation
 15 steps, the proposed semi-supervised clustering will converge to a (at least
 16 local) minimum in each interactive iteration.

17 After each interactive iteration, new constraints are given to the system.
 18 These new constraints might be in contradiction with some of the ones pre-
 19 viously deduced by the system from the earlier user interactive iterations.
 20 For this reason and also for computational time matters, our system omits

1 at each step some of the constraints deduced at earlier steps. Therefore, the
2 objective function J_{obj} may be different between different interactive itera-
3 tions. And the convergence of the interactive semi-supervised model is thus
4 not guaranteed. But we can verify the convergence of the model, practically,
5 by determining, at the end of all interactive iterations, the global objective
6 function which considers all feedback given by the user in all interactive it-
7 erations and then by verifying if this global objective function has improved
8 or not after different interactive steps. This is a part of our current work.

9 *3.3. Interactive interface*

10 In order to allow the user to view the clustering results and to provide
11 feedback to the systems, we implement an interactive interface as shown in
12 Figure 1.

13 The rectangle at the bottom right corner of Figure 1 is the principal
14 plane representing all presented clusters by their prototype images. In our
15 system, the maximum number of cluster prototypes presented to the user on
16 the principal plane is fixed at 30. The prototype image of each cluster is the
17 most representative image of that cluster chosen as follows. In our model,
18 we use the internal measure Silhouette-Width (SW) (Rousseeuw, 1987) to
19 estimate the quality of each image in a cluster. The higher the SW value
20 of an image in a cluster, the more representative this image is for the clus-
21 ter. The prototype image of a cluster is thus the image with the highest
22 SW value in the cluster. Any other internal measure could be used instead.
23 The position of the prototype image of each cluster in the principal plane
24 represents the position of the corresponding cluster center. It means that, if
25 two cluster centers are close (or distant) in the n-dimensional feature space,

1 their prototype images are close (or distant) in the 2D principal plane. For
2 representing the cluster centers which are n-dimensional vectors in 2D plane,
3 we use Principal Component Analysis (PCA) (Pearson, 1901); the principal
4 plane consists of the two principal axes associated with the highest eigen-
5 values. The importance of an axis is represented by its inertia (the sum of
6 the squared elements of this axis (Abdi and Williams, 2010)) or by the per-
7 centage of its inertia in the total inertia of all axes. In general, if the two
8 principal axes explain (cumulatively) greater or equal to 80% of the total
9 inertia, the PCA approach could lead to a nice 2D-representation of the pro-
10 totype images. In our case, the accumulated inertia explained by the two
11 first principal axes is about 65% for the Wang and PascalVoc2006 databases
12 and about 20% for the Caltech101 and Corel30k image databases. As only a
13 maximum of 30 clusters (and therefore 30 prototype images) can be shown
14 to the user in an interactive iteration, a not very nice 2D-representation of
15 prototype images does not influence on the results as long as the user can
16 distinguish between the prototype images and have a rough idea of the dis-
17 tances between the clusters. When there are some prototype images which
18 overlap each other, a slight modification of the PCA components can help to
19 separate these images.

20 By clicking on a prototype image in the principal plane, the user can view
21 the corresponding cluster. In Figure 1, each cluster selected by the user is
22 represented by a circle:

- 23 • The prototype image of this cluster is located at the center of the circle.
- 24 • The 10 most representative images (images with the highest SW val-
25 ues), which have not received feedback from the user in the previous

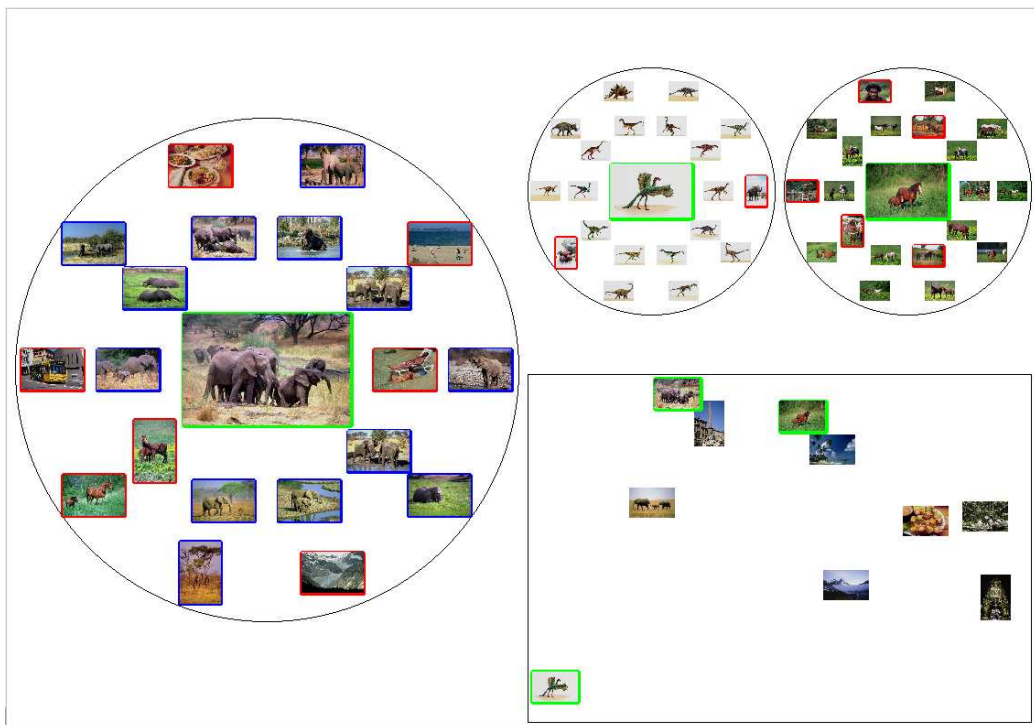


Figure 1: 2D interactive interface. The rectangle at the bottom right corner represents the principal plane consisting of the two first principal axes (obtained by PCA) of the prototype images of all clusters. Each circle represents the details of a particular cluster selected by the user.

1 iterations, are located in the first circle of images around the prototype
2 image, near the center.

3 • The 10 least representative images (images with the smallest SW val-
4 ues), which have not received feedback from the user in the previous
5 iterations, are located in the second circle of images around the proto-
6 type image, close to the cluster border.

7 By showing, for each iteration, the images which have not received user feed-
8 back in previous iterations, we wish to obtain feedback for different images.

9 The user can specify positive feedback and negative feedback (images in
10 Figure 1 with blue and red borders respectively) for each cluster. The user
11 can also change the cluster assignment of a given image by dragging and
12 dropping the image from the original cluster to the new cluster. When an
13 image is changed from cluster A to cluster B, it is considered as negative
14 feedback for cluster A and positive feedback for cluster B. Therefore, after
15 each interactive iteration, the process returns a positive image list and a
16 negative image list for each cluster with which the user has interacted.

17 *3.4. Pairwise constraint deduction*

18 In each interactive iteration, user feedback is in the form of positive and
19 negative images, while the supervised input information of the proposed semi-
20 supervised clustering method are pairwise constraints between CF entries.
21 Therefore, we have to deduce the pairwise constraints between CF entries
22 from the user feedback.

23 At each interactive iteration and for each interacted cluster, all positive
24 images should be in this cluster while negative images should move to another

1 cluster. We consider that each image in the positive set is linked to each image
 2 in the negative set by a cannot-link, while all images in the positive set are
 3 linked by must-links. If we assume that all feedback is coherent between
 4 different interactive iterations, we try to group images, which should be in
 5 the same cluster according to the user feedback of all interactive iterations,
 6 in a group called *neighborhood*. We define:

- 7 • $Np = \{Np_i\}$ is the neighborhood list, each neighborhood $Np_i = \{x_j\}$
 8 including a list of images which should be in a same cluster.
- 9 • $CannotNp = \{cannotNp_i\}$, each element $cannotNp_i = \{n_j\}$ including
 10 labels of the neighborhoods which should not be in the same cluster as
 11 Np_i . Two neighborhoods Np_i and Np_j are called cannot-link neigh-
 12 borhoods if there is at least one cannot-link between a point of Np_i
 13 and a point of Np_j .

14 After receiving the list of feedback in the current iteration, the lists Np and
 15 $CannotNp$ are updated as follows:

- 16 1. Update based on positive feedback: For each cluster μ_h which receives
 17 interaction from the user:
 - 18 (a) Initialize $n_h \leftarrow -1$, n_h indicates the neighborhood including pos-
 19 itive images of the cluster μ_h .
 - 20 (b) If all positive images of μ_h are not included in any neighborhood
 21 \rightarrow create a new neighborhood for these positive images and assign
 22 n_h as the index of this neighborhood.

- 1 (c) If some positive images of μ_h are already included in one or mul-
 2 tiple neighborhoods \rightarrow merge these neighborhoods (in the case of
 3 multiple neighborhoods) into one single neighborhood, insert the
 4 other positive images which are not included in any neighborhood
 5 to this neighborhood and update n_h as the index of this neigh-
 6 borhood. Also update the set *CannotNp* to signify that neighbor-
 7 hoods that had cannot-link with one of the neighborhoods which
 8 has merged, now have cannot-link with the new neighborhood.
- 9 2. Update based on negative feedback: For each negative image x_j of each
 10 cluster μ_h which receives interaction from the user:
- 11 (a) If x_j is not included in any neighborhood \rightarrow create a new neigh-
 12 borhood for x_j .
- 13 (b) If x_j is already included in the neighborhood Np_{n_j} , and Np_{n_h}
 14 is the neighborhood corresponding to the positive images of the
 15 cluster μ_h , update the corresponding *cannotNp_{n_j}* and *cannotNp_{n_h}*
 16 to signify that Np_{n_j} and Np_{n_h} have cannot-link.

17 As we assume that the user feedback is coherent among different inter-
 18 active iterations, all images in a same neighborhood should be in a same
 19 cluster and images of cannot-link neighborhoods should be in different clus-
 20 ters. There may be cannot-link images belonging to the same CF_i . There
 21 may also be simultaneous must-link and cannot-link between images of CF_i
 22 and images of CF_j . In such cases, these CF entries should be split into purer
 23 CF entries. To do so, we define a *seed* of an entry CF_i as a subset of images
 24 of CF_i so that the images of this *seed* are included in a same neighborhood.

1 Therefore, an entry CF_i may contain some *seeds* corresponding to different
2 neighborhoods and other images which are not included in any other neigh-
3 borhood. Cannot-link may or may not exist between seeds of a CF entry.
4 With each CF entry that should be split, we present the user with each pair
5 of seeds, which do not have cannot-link between them, to demand more in-
6 formation (for each seed, the image which is closest to the center of the seed
7 is presented):

- 8 • If the user indicates that there is must-link between these two seeds,
9 these seeds and also their corresponding neighborhoods are merged.
- 10 • If the user indicates that there is cannot-link between these two seeds,
11 update the corresponding *cannotCF* lists specifying that their two cor-
12 responding neighborhoods have cannot-link between them.

13 An entry CF_i is split as follows: if CF_i has p seeds, it should be split into
14 p different CF entries; each new CF entry contains all points of a seed; ev-
15 ery other point of CF_i which is not included in any seed is assigned to the
16 CF entry corresponding to the closest seed. By splitting the necessary CF
17 entries into purer CF entries, we can eliminate the case where cannot-link
18 exists between images of a same CF or where must-link and cannot-link ex-
19 ist simultaneously between images of two different CF entries. Subsequently,
20 pairwise constraints between CF entries can be deduced based on pairwise
21 constraints between images as follows: if there is must-link (or respectively
22 cannot-link) between two images of two CF entries, a must-link (or respec-
23 tively cannot-link) is created between these two CF entries.

24 Concerning pairwise constraints between images, a simple and complete

1 way to deduce them is to create must-link between each pair of images of a
2 same neighborhood, and to create, for each pair of cannot-link neighborhoods
3 (Np_i, Np_j) , cannot-link between each image of Np_i and each image of Np_j .
4 By deducing pairwise constraints between images in this way, the number
5 of constraints between images can be very high, and therefore the number
6 of constraints between CF entries could also be very high. The processing
7 time of the semi-supervised clustering in the next phase could thus be very
8 high due to the high number of constraints. There are different strategies for
9 deducing pairwise constraints between images that could reduce the number
10 of constraints and also the processing time. One of them is presented in
11 Figure 2 and others are described and tested in Section 4. In Figure 2,
12 must-links are created between positive images of each cluster while cannot-
13 link are created between positive and negative images of each cluster (note
14 the displacement feedback corresponding to a negative image of the source
15 cluster and a positive image of the destination cluster).

16 **4. Experiments**

17 In this section, we present some experimental results of our interactive
18 semi-supervised clustering model. We also, experimentally, compare our
19 semi-supervised clustering model with the semi-supervised HMRF-kmeans.
20 When using the semi-supervised HMRF-kmeans in the re-clustering phase,
21 the initial unsupervised clustering is k-means.

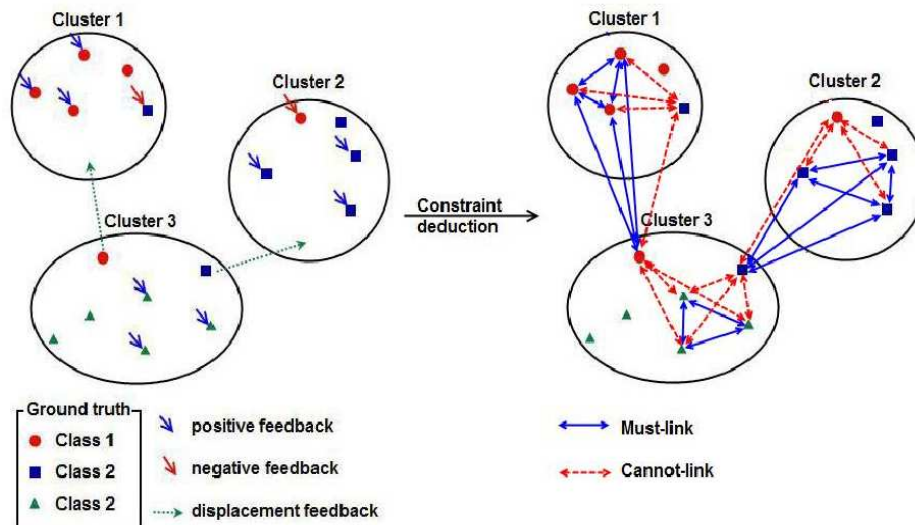


Figure 2: Example of pairwise constraint deduction between images from the user feedback.

1 4.1. Experimental protocol

2 In order to analyze the performance of our interactive semi-supervised
 3 clustering model, we use different image databases (Wang¹ (1000 images di-
 4 vided into 10 classes), PascalVoc2006² (5304 images divided into 10 classes),
 5 Caltech101³ (9143 images divided into 101 classes)). Note that in our ex-
 6 periments we use the same number of clusters as the number of classes in
 7 the ground truth. As presented in Section 3.3, the cluster prototype images
 8 are shown to the user on the principal plane; users can choose to view and
 9 interact with any cluster in which they are interested. For databases which

¹<http://wang.ist.psu.edu/docs/related/>

²<http://pascallin.ecs.soton.ac.uk/challenges/VOC/>

³<http://www.vision.caltech.edu/ImageDatasets/Caltech101/>

1 have a small number of classes, such as Wang and PascalVoc2006, all proto-
2 type images can be shown on the principal plane. For databases which have
3 a large number of classes, such as Caltech101, only a part of the prototype
4 images can be shown for visualization. In our system, the maximum number
5 of cluster prototypes shown to the user in each iteration is fixed at 30. We use
6 two simple strategies for choosing clusters to be shown for each iteration: 30
7 clusters chosen randomly or iteratively chosen pairs of closest clusters until
8 there are 30 clusters.

9 The external measures compare the clustering results with the ground
10 truth, thus they are compatible for estimating the quality of the interactive
11 clustering involving user interaction. As different external measures analyze
12 the clustering results in a similar way (see Lai et al. (2012a)), we use, in this
13 paper, the external measure V-measure (Rosenberg and Hirschberg, 2007).
14 The greater the V-measure values are, the better the results (compared to
15 the ground-truth).

16 Concerning feature descriptors, we implement the local descriptor rgSIFT
17 (van de Sande et al., 2008), an extension for color image of the SIFT descrip-
18 tor (Lowe, 2004), that today is widely used for its high performance. The
19 SIFT descriptor detects interest points from an image and describes the local
20 neighborhood around each interest point by a 128-dimensional histogram of
21 local gradient directions of image intensities. The rgSIFT descriptor of each
22 interest point is computed as the concatenation of the SIFT descriptors calcu-
23 lated for the r and g components of the normalized RGB color space (van de
24 Sande et al., 2008) and the SIFT descriptor in the intensity channel, resulting
25 in a 3×128 -dimensional vector. The “Bag of words” (Sivic and Zisserman,

1 2003) approach is chosen to group local features of each image into a single
2 vector. It consists in two steps. Firstly, K-means clustering is used to group
3 local features of all images in the database according to a number *dictSize*
4 of clusters. We then generate a dictionary containing *dictSize* visual words
5 which are the centroids of these clusters. The feature vector of each image
6 is a *dictSize* dimension histogram representing the frequency of occurrence
7 of the visual words in the dictionary, by replacing each local descriptor of
8 the image by the nearest visual word. Our experiments in (Lai et al., 2012a)
9 show that local descriptors are better than global descriptors regarding the
10 external measures and the value *dictSize* = 200 is a good trade-off between
11 the size of the feature vector and the performance. Therefore, in our experi-
12 ments, we use the rgSIFT descriptor together with a visual word dictionary
13 of size 200.

14 In order to undertake the interactive tests automatically, we implement a
15 software agent, later referred to as “user agent” that simulates the behavior
16 of the human user when interacting with the system (assuming that the agent
17 knows all the ground truth containing the class label for each image). At each
18 interactive iteration, clustering results are returned to the user agent by the
19 system; the agent simulates the behavior of the user giving feedback to the
20 system. For simulating the user behavior, we suggest some rules:

- 21 • At each interactive iteration, the user agent interacts with a fixed num-
22 ber of c clusters.
- 23 • The user agent uses two strategies for choosing clusters: randomly
24 chosen c clusters, or iteratively chosen pairs of closest clusters until
25 there are c clusters.

- 1 • The user agent determines the image class (in the ground truth) cor-
2 responding to each cluster by the most represented class among the 21
3 presented images of the cluster. The number of images of this class
4 in the cluster must be greater than a threshold *MinImages*. If this is
5 not the case, this cluster can be considered as a noise cluster. In our
6 experiments, $MinImages = 5$ for databases having a small number
7 of classes (Wang, PascalVoc2006), and $MinImages = 2$ for databases
8 having a large number of classes (Caltech101).

- 9 • When several clusters (among chosen clusters) correspond to a same
10 class, the cluster in which the images of this class are the most numerous
11 (among the 21 shown images of the cluster) is chosen as the principal
12 cluster of this class. The classes of the other clusters are redefined as
13 usual, but neutralize the images from this class.

- 14 • In each chosen cluster, all images, where the result of the algorithm
15 corresponds to the ground truth, are labeled as positive samples of
16 this cluster, while the others are negative samples of this cluster. All
17 negative samples are moved to the cluster (among chosen clusters) cor-
18 responding to their class in the ground truth.

19 As presented in Section 3.4, we have to deduce pairwise constraints be-
20 tween images based on user feedback in each iteration and also on the neigh-
21 borhood information. User feedback is in the form of positive and negative
22 images of each cluster (the image which is displaced from one cluster to an-
23 other cluster is considered as a negative image of the source cluster and a
24 positive image of the destination cluster). The neighborhood information is

1 in the form of the lists $Np = \{Np_i\}$ and $CannotNp = \{cannotNp_i\}$, where
2 each neighborhood Np_i contains images which should be in a same cluster
3 and $cannotNp_i$ identifies the list of neighborhoods having cannot-link with
4 Np_i . Neighborhood information is deduced from user feedback during all
5 interactive iterations, as presented in Section 3.4. Pairwise constraints be-
6 tween images will be used directly for the semi-supervised HMRF-Kmeans,
7 while they have to be deduced into pairwise constraints between CF entries
8 (see Section 3.4) to be used by our proposed semi-supervised clustering. We
9 divide pairwise constraints between images into two kinds: *user constraints*
10 and *deduced constraints*. *User constraints* are created directly, based on user
11 feedback in each iteration, while *deduced constraints* are created by deduction
12 rules. For instance, in the first iteration, the user marks x_1, x_2 as positive
13 images and x_3 as a negative image of cluster μ_i ; while in the second iter-
14 ation, he marks x_1 and x_4 as positive images of cluster μ_j . The created
15 user constraints are: must-link between positive images in the first iteration
16 (x_1, x_2) , must-link between positive images in the second iteration (x_1, x_4) ,
17 and cannot-links between positive and negative images in the first iteration
18 $(x_1, x_3), (x_2, x_3)$. As there are must-links $(x_1, x_2), (x_1, x_4)$, there is also a de-
19 duced must-link (x_2, x_4) . In addition deduced cannot-link (x_3, x_4) is created,
20 based on the must-link (x_1, x_4) and the cannot-link (x_1, x_3) . We can see that
21 deduced constraints can be created based on neighborhood information. In
22 our experiments, we use different strategies for deducing pairwise constraints
23 between images. These strategies are detailed in Table 1.

Table 1: Different strategies for deducing pairwise constraints between images based on user feedback and on neighborhood information.

N ^o	Take into account	Details
1	<ul style="list-style-type: none"> • All user constraints of all interactive iterations. • All deduced constraints of all interactive iterations. 	<p>All constraints are created based on the neighborhood information:</p> <ul style="list-style-type: none"> • Must-link between each pair of images of each neighborhood. • Cannot-link between each image of each neighborhood $Np_i \in Np$ and each image of each neighborhood having cannot-link with Np_i (listed in <i>cannotNp_i</i>).
2	<ul style="list-style-type: none"> • All user constraints of all interactive iterations. • None of deduced constraints. 	<p>In each iteration, all possible user constraints are created:</p> <ul style="list-style-type: none"> • Must-link between each pair of positive images of each cluster. • Cannot-link between each pair of a positive image and a negative image of a same cluster.
3	<ul style="list-style-type: none"> • All user constraints of all interactive iterations. • All deduced constraints in the current iteration (deduced constraints in the previous iterations are eliminated). 	<ul style="list-style-type: none"> • In each iteration, all possible user constraints are created as in Strategy 2. • Deduced constraints in the current iteration are created while updating the neighborhoods as follows: <ul style="list-style-type: none"> – If there is a must-link (or cannot-link) (x_i, x_j), $x_j \in Np_m$, deduced must-links (or cannot-links) (x_i, x_l), $\forall x_l \in Np_m$ are created. – If there is a must-link (or cannot-link) (x_i, x_j), $x_i \in Np_m$, $x_j \in Np_n$, deduced must-links (or cannot-links) (x_k, x_l), $\forall x_k \in Np_m$, $\forall x_l \in Np_n$ are created.
Continued on next page		

Table 1 – continued from previous page

N°	Take into account	Details
4	<ul style="list-style-type: none"> • User constraints between images and cluster centers of <i>all</i> interactive iterations. • Deduced constraints between images and cluster centers in the <i>current</i> iteration (deduced constraints in the previous iterations are eliminated). 	<p>In each iteration, the positive image having the best internal measure (SW) value among all positive images of each cluster is the center of this cluster.</p> <ul style="list-style-type: none"> • Must-link/cannot-link user constraints are created in each iteration between each positive/negative image and the corresponding cluster center. • Deduced constraints in the current iteration are created while updating the neighborhoods as follows: <ul style="list-style-type: none"> – If x_i and x_j must be in the same (or different) clusters (based on user feedback), $x_j \in Np_m$, deduced must-links (or cannot-links) are created between x_i and each center image of Np_m. – If x_i and x_j must be in the same (or different) clusters (based on user feedback), $x_i \in Np_m$, $x_j \in Np_n$, deduced must-links (or cannot-links) are created between x_i and each center image of Np_n and between x_j and each center image of Np_m.
5	<ul style="list-style-type: none"> • User constraints (must-links between the most distant images and cannot-links between the closest images) of <i>all</i> iterations. • Deduced constraints (must-links between the most distant images and cannot-links between the closest images) of <i>all</i> iterations. 	<ul style="list-style-type: none"> • User constraints are created for each cluster in each iteration as follows: must-links are successively created between two positive images (at least one of them is not selected by any must-link) that have the longest distance until all positive images of the cluster are connected by these must-links; cannot-links are created between each negative image and the nearest positive image of the cluster. • Deduced constraints are created in each iteration as follows: must-links for each neighborhood are successively created between two images that have the longest distance until all images of this neighborhood are connected by these must-links; cannot-links are deduced, for each pair of cannot-link neighborhoods (Np_i, Np_j), between each image of Np_i and the nearest image of Np_j and between each image of Np_j and the nearest image of Np_i.
6	<p>Same idea as in strategy 5, but the size of the neighborhoods is considered while creating deduced cannot-links.</p>	<p>User constraints and deduced must-link constraints are created as in Strategy 5. For each pair of cannot-link neighborhoods, deduced cannot-links are only created between each image of the neighborhood that has the least number of images and the nearest image of the neighborhood that has the most images.</p>

1 *4.2. Experimental results*

2 *4.2.1. Analysis of different strategies for deducing pairwise constraints be-*
3 *tween images*

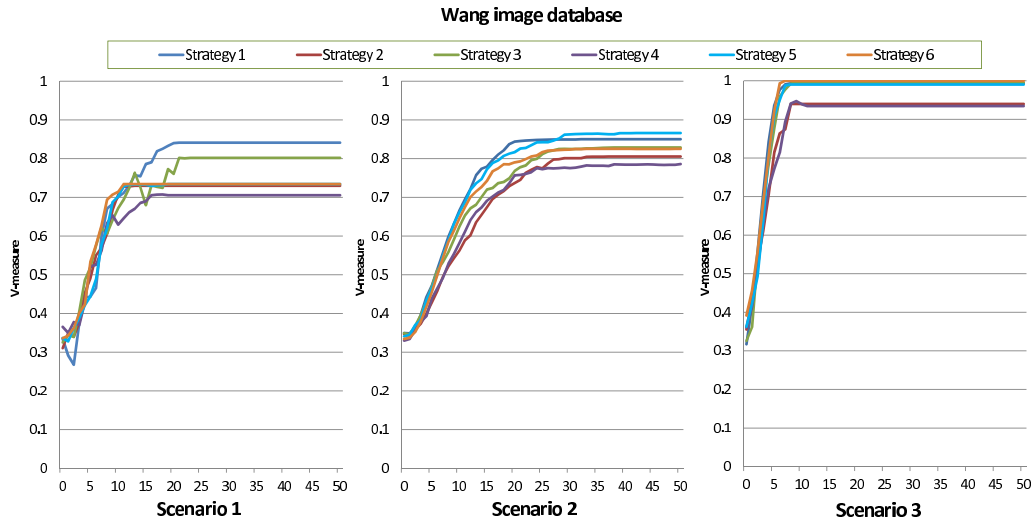
4 The first set of experiments aims at evaluating the performance of our
5 interactive semi-supervised clustering model using different strategies for de-
6 ducing pairwise constraints between images. Note that constraints between
7 CF entries should be deduced from constraints between images, before being
8 used in the re-clustering phase. We use the Wang and the PascalVoc2006
9 image databases for these experiments. For these two databases, we propose
10 three test scenarios (note that c specifies the number of clusters which are
11 chosen for interacting in each iteration):

- 12 • Scenario 1: $c = 5$ closest clusters are chosen.
- 13 • Scenario 2: $c = 5$ clusters are randomly chosen.
- 14 • Scenario 3: $c = 10$, all cluster are chosen (Wang and PascalVoc2006
15 both have 10 clusters).

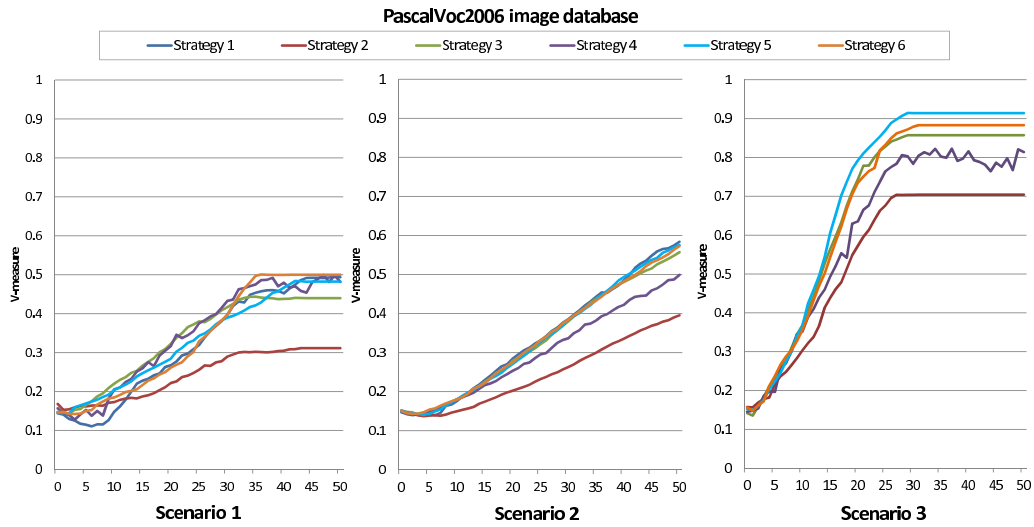
16 Note that our experiments are carried out automatically, *i.e.* the feedback is
17 given by a software agent simulating the behaviors of the human user when
18 interacting with the system. In fact, the human user can give feedback by
19 clicking for specifying the positive and/or negative images of each cluster or
20 by dragging and dropping the image from a cluster to another cluster. For
21 each cluster selected by the user, only 21 images of this cluster are displayed
22 (see Figure 1). Therefore, for interacting with 5 clusters (scenarios 1, 2) or
23 10 clusters (scenario 3), the user has to realize respectively a maximum of
24 105 or 210 mouse clicks in each interactive iteration. These upper bounds do

1 not depend on neither the size of the database nor the pairwise constraint
2 deduction strategy, and in practice the number of clicks that the user has
3 to provide is far lower. However, the number of deduced constraints may be
4 much greater than the user’s clicks (and this number depends on the database
5 size and on the pairwise constraint deduction strategy). When applying the
6 interactive semi-supervised clustering model in the indexing phase, the user
7 is generally required to provide as much feedback as possible for having a
8 good indexing structure which could lead to better results in the further
9 retrieval phase. Therefore, in the case of the indexing phase, the proposed
10 number of clicks seems tractable.

11 Figures 3a and 3b show, respectively, the results during 50 interactive
12 iterations of our proposed interactive semi-supervised clustering model on
13 the Wang and PascalVoc2006 image databases, with the three proposed sce-
14 narios. The results are shown according to 6 strategies for deducing pairwise
15 constraints presented in Table 1. The vertical axis specifies the V-measure
16 values, while the horizontal axis specifies the number of iterations. Note
17 that with each selected cluster, the user agent gives all possible feedback.
18 Therefore, for each scenario, the numbers of user feedback are equivalent
19 between different iterations and between different strategies. As in scenario
20 2, clusters are randomly chosen, we realize this scenario 10 times for each
21 database. The curves of the scenario 2 shown in Figures 3a and 3b represent
22 the mean values of the V-measure over these 10 executions at each iteration.
23 The average standard deviation of each strategy after 50 iterations is pre-
24 sented in Table 2. The corresponding execution time for these experiments
25 is presented in Table 3 (note that for the scenario 2, the average execution



(a) Results on the Wang image database



(b) Results on the PascalVoc2006 image database

Figure 3: Results of our proposed interactive semi-supervised clustering model during 50 interactive iterations on the Wang and PascalVoc2006 image databases, using 6 strategies for deducing pairwise constraints. The horizontal axis specifies the number of iterations.

Table 2: Average standard deviation of 10 executions of the scenario 2 after 50 interactive iterations corresponding to the experiments of our proposed interactive semi-supervised clustering model shown in Figures 3a and 3b

	Average standard deviation	
	Wang database	PascalVoc2006 database
Strategy 1	0.033	0.022
Strategy 2	0.044	0.017
Strategy 3	0.045	0.025
Strategy 4	0.047	0.022
Strategy 5	0.036	0.024
Strategy 6	0.044	0.026

1 times of 10 executions are shown). The experiments are executed using a
 2 normal PC with 2GB of RAM.

3 We can see that the clustering results progress, in general, after each in-
 4 teractive iteration, in which the system re-clusters the dataset by considering
 5 the constraints deduced from accumulated user feedback. In most cases, the
 6 clustering results converge after only a few iterations. This may be due to the
 7 fact that no new knowledge is provided. Moreover, we can easily see that the
 8 clustering results are better and converge more quickly when the number of
 9 chosen clusters (and therefore the number of constraints) in each interactive
 10 iteration is higher (scenario 3 gives better results and converges more quickly
 11 than scenarios 1 and 2). In addition, for both image databases, scenario 2,
 12 in which clusters are randomly chosen for interacting, gives better results
 13 than scenario 1, in which the closest clusters are chosen. When selecting the

Table 3: Processing time after 50 interactive iterations of the experiments of our proposed interactive semi-supervised clustering model shown in Figures 3a and 3b.

	Wang database		
	Scenario 1	Scenario 2	Scenario 3
Strategy 1	1h58'	2h24'	1h41'
Strategy 2	9'	12'	10'
Strategy 3	31'	19'	47'
Strategy 4	8'	9'	8'
Strategy 5	8'	9'	9'
Strategy 6	6'	8'	8'
	PascalVoc2006 database		
	Scenario 1	Scenario 2	Scenario 3
Strategy 1	16d12h	14d11h	
Strategy 2	2h55'	4h02'	5h6'
Strategy 3	3h23'	6h39'	6h22'
Strategy 4	1h9'	1h33'	2h17'
Strategy 5	3h33'	4h42'	3h10'
Strategy 6	1h3'	1h21'	2h

1 closest clusters there may be only several clusters that always receive user
2 feedback; thus the constraint information is less than when all the clusters
3 could receive user feedback when we randomly select the clusters.

4 As regards different strategies for deducing pairwise constraints, we can
5 see that for each database, the average standard deviations over 10 executions
6 of the scenario 2 are similar for all scenarios. Therefore, we can compare
7 different strategies based on the mean values shown on Figures 3a and 3b.
8 We can see that:

- 9 • Strategy 1 shows, in general, very good performance but the processing
10 time is huge because it uses all possible user constraints and deduced
11 constraints created during all iterations.
- 12 • Strategy 2, the only strategy uniquely using user constraints, generally
13 gives the worst results; thus deduced constraints are needed for better
14 performance. Its processing time is also high due to the large number
15 of user constraints.
- 16 • Strategy 3 shows good or very good performance but some oscillations
17 exist between different iterations because, when overlooking previously
18 deduced constraints, some important constraints may be omitted. Its
19 processing time is high.
- 20 • Strategy 4 gives better results than strategy 2, but the results are unsta-
21 ble because this strategy also overlooks previously deduced constraints.
22 It has good execution time while reducing the number of constraints.
- 23 • Strategy 5 generally gives good or very good results by keeping im-
24 portant constraints (must-links between the most distant images and

Table 4: Processing time after 50 interactive iterations corresponding to the experiments presented in Figures 4a and 4b of the proposed semi-supervised clustering and of the semi-supervised HMRF-kmeans. Strategy 6 in Table 1 for deducing pairwise constraints is used.

	Wang database		
	Scenario 1	Scenario 2	Scenario 3
Proposed semi-supervised clustering	6'	8'	8'
HMRF-kmeans	7'	11'	10'
	PascalVoc2006 database		
	Scenario 1	Scenario 2	Scenario 3
Proposed semi-supervised clustering	1h3'	1h21'	2h'
HMRF-kmeans	2h16'	3h10'	2h49'

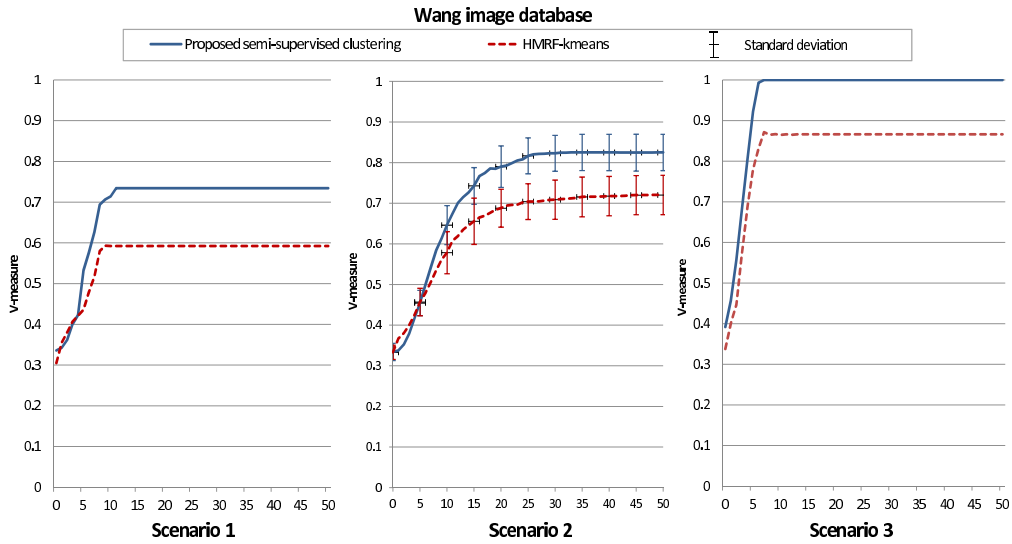
1 cannot-links between the closest images), but its processing time is still
2 high.

- 3 • Strategy 6, by reducing the deduced cannot-link constraints from strat-
4 egy 5, gives in general very good results in low execution time.

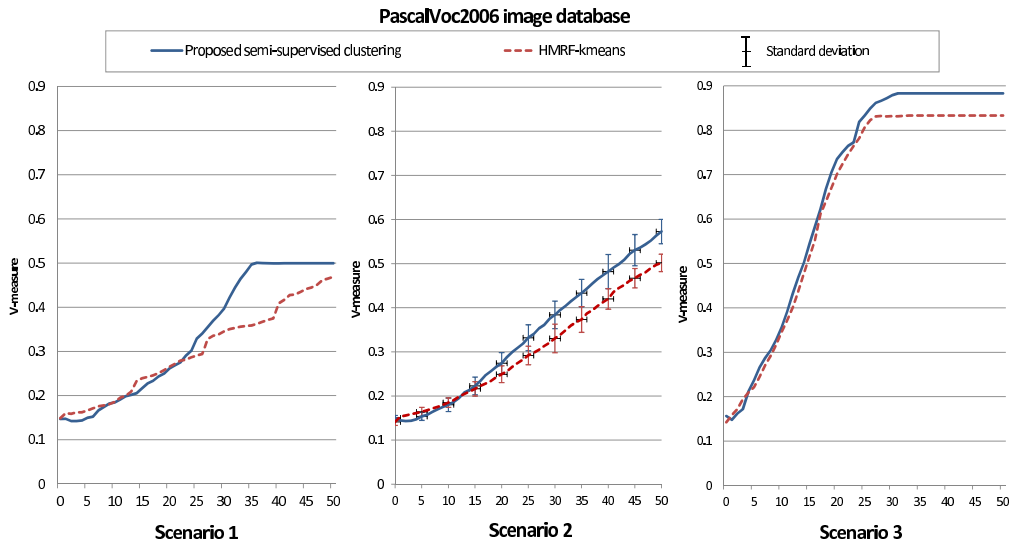
5 We can conclude, from this analysis, that strategy 6 shows the best trade-
6 off between performance and processing time. This strategy will be used in
7 further experiments.

8 *4.2.2. Comparison of the proposed semi-supervised clustering model and the* 9 *semi-supervised HMRF-kmeans*

10 Figures 4a and 4b represent, respectively, the clustering results for 50
11 interactive iterations on the Wang and the PascalVoc2006 image databases



(a) Comparison results on the Wang image database



(b) Comparison results on the PascalVoc2006 image database

Figure 4: Comparison of the proposed semi-supervised clustering and the semi-supervised HMRF-kmeans with 50 interactive iterations using Strategy 6 in Table 1 for deducing the pairwise constraints between images. The horizontal axis represents the number of iterations.

1 when using our proposed semi-supervised clustering and the semi-supervised
2 HMRF-kmeans in the re-clustering phase. The three scenarios described
3 in Section 4.2.1 and strategy 6, for deducing pairwise constraints between
4 images, are used. Note that the results of scenario 2 represent the mean
5 values and also the standard deviations over 10 executions at each itera-
6 tion. The corresponding processing time is presented in Table 4. We can
7 see that in all scenarios, our proposed method gives better results, in lower
8 processing time than the HMRF-kmeans. While the pairwise constraints
9 between images are directly used by the HMRF-kmeans, they are deduced
10 in pairwise constraints between CF entries for being used by our proposed
11 semi-supervised clustering. A CF entry groups a list of similar images, thus
12 many pairwise constraints between images can be represented by only one
13 pairwise constraints between CF entries. Therefore, with a same set of user
14 feedback, the number of pairwise constraints between images is generally
15 greater than the number of the pairwise constraints between CF entries.
16 Thus the processing time of the HMRF-kmeans is much higher than the pro-
17 cessing time of our proposed method. Moreover, when a pairwise constraint
18 (CF_i, CF_j) is deduced from the pairwise constraint of the corresponding im-
19 ages $(x_k, x_l), x_k \in CF_i, x_l \in CF_j$, the constraint (CF_i, CF_j) forces the group-
20 ing or separating of not only the two images x_i and x_j but also the other im-
21 ages included in CF_i and CF_j . And therefore, the clustering results given by
22 our proposed method are better than the ones given by the HMRF-kmeans.
23 Moreover, similar to the experiments presented in Section 4.2.1, the scenario
24 2 in which the clusters are randomly chosen for interacting gives better re-
25 sults than the scenario 1 in which the closest clusters are chosen. In the

Table 5: Processing time after 50 interactive iterations corresponding to the experiments on the Caltech101 image database in Figure 5 for the proposed semi-supervised clustering and for the semi-supervised HMRF-kmeans. Strategy 6 in Table 1 for deducing pairwise constraints is used.

	Proposed semi-supervised clustering	HMRF-kmeans
Scenario 4	13h26'	48h33'
Scenario 5	8h4'	33h45'
Scenario 6	33h34'	157h26'
Scenario 7	50h12'	101h11'

1 following experiments on the Caltech101 image database, we present only
 2 the clustering results when the clusters are randomly chosen.

3 As the Caltech101 database has a large number of classes (101 classes),
 4 we do not show all clusters to the user on the principal plane but only a
 5 small number of clusters (we fix the maximum number of cluster that could
 6 be shown on the principal plane to 30). There are two strategies for choosing
 7 clusters to be shown on the principal plane: either clusters are randomly
 8 chosen or the closest clusters are chosen. The user agent randomly chooses,
 9 among shown clusters, c clusters for interacting. We use 4 scenarios for the
 10 experiments on the Caltech101 image database:

- 11 • Scenario 4: the closest clusters are chosen to be shown to the user, $c=5$
 12 clusters are chosen by the user agent for interacting.
- 13 • Scenario 5: clusters are randomly chosen to be shown to the user, $c=5$
 14 clusters are chosen by the user agent for interacting.

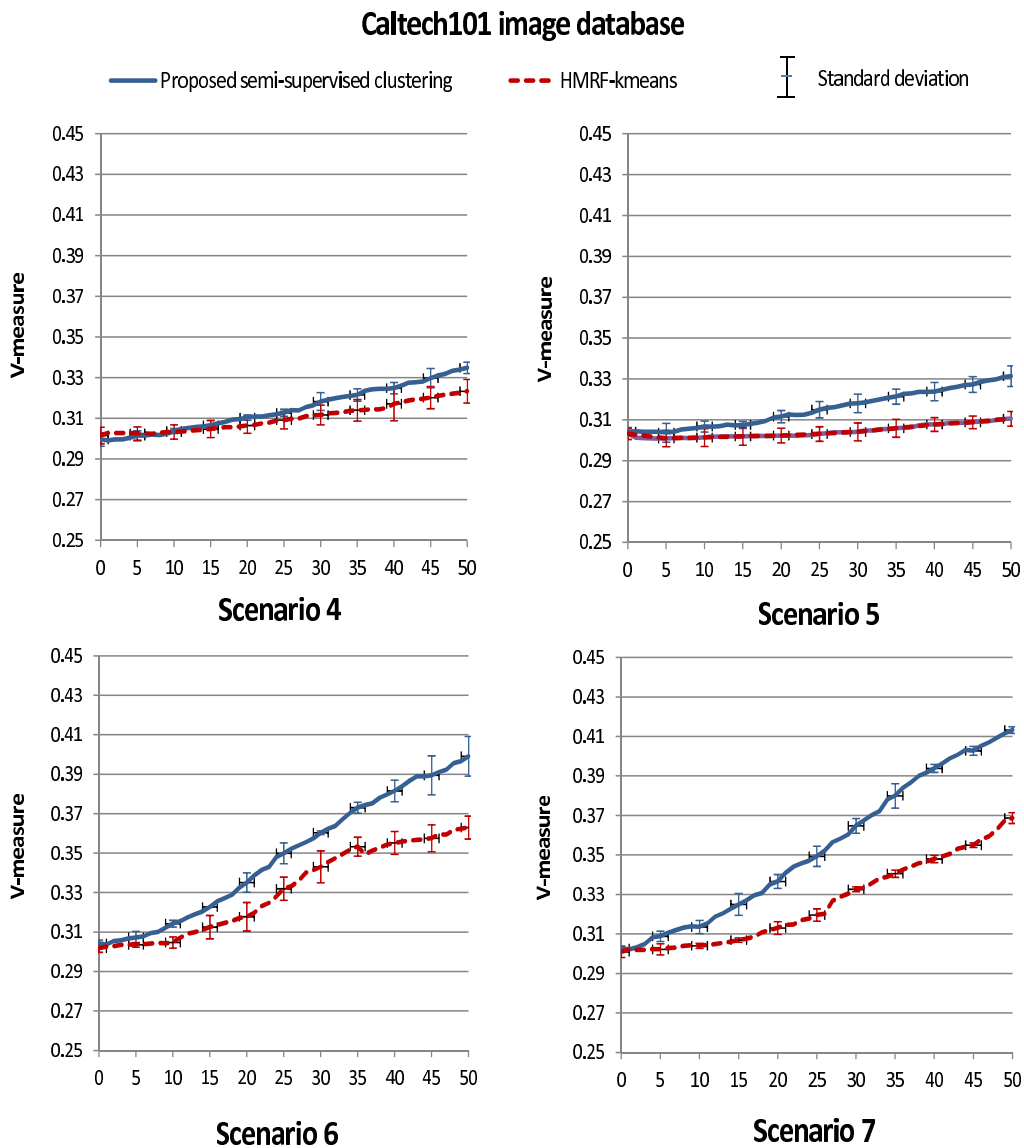


Figure 5: Comparison of the proposed semi-supervised clustering and the semi-supervised HMRF-kmeans on the Caltech101 image database for 50 interactive iterations. The strategy 6 in Table 1 for deducing pairwise constraints are used. The horizontal axis represents the number of iterations.

- 1 • Scenario 6: the closest clusters are chosen to be shown to the user,
2 $c=10$ clusters are chosen by the user agent for interacting.
- 3 • Scenario 7: clusters are randomly chosen to be shown to the user, $c=10$
4 clusters are chosen by the user agent for interacting.

5 Figure 5 compares our proposed semi-supervised clustering and the HMRF-
6 kmeans during 50 interactive iterations on the Caltech101 image database.
7 The corresponding processing time is presented in Table 5. As in all these
8 four scenarios, the clusters are randomly chosen for interacting, we realize
9 each scenario 5 times and present in Figure 5 the mean values and also the
10 standard deviations over 5 executions. The results shows that our proposed
11 semi-supervised clustering outperforms the HMRF-kmeans in all four sce-
12 narios. Moreover, the clustering results are also better when the number of
13 feedback for each iteration is high (scenarios 6 and 7 give better results than
14 scenarios 4 and 5).

15 5. Conclusion

16 A new interactive semi-supervised clustering model for indexing image
17 databases is presented in this article. After receiving user feedback for each
18 interactive iteration, the proposed semi-supervised clustering re-organizes the
19 dataset by considering the pairwise constraints between CF entries deduced
20 from the user feedback. We present an interactive interface allowing the
21 user to view, and to provide feedback. Experimental analysis, using a soft-
22 ware user agent for simulating human user behavior, shows that our model
23 improves the clustering results at each interactive iteration. [Note that our](#)

1 experimental scenarios are realistic, they can be realized by a real user as the
2 number of clicks required is tractable. The experiments on different image
3 databases (Wang, PascalVoc2006, Caltech101), presented in this paper, also
4 show that our semi-supervised clustering outperforms the semi-supervised
5 HMRF-kmeans (Basu et al., 2004) in both performance and processing time.

6 Moreover, we propose and compare, experimentally, different strategies
7 for deducing pairwise constraints from the user feedback accumulated from all
8 interactive iterations. The experimental results show that strategy 6 in Table
9 1, which keeps only the most important constraints (must-links between the
10 most distant images and cannot-links between the closest images), provides
11 the best trade-off between the performance and the processing time. Strategy
12 6 is therefore the most suitable, in our context involving the user in the
13 indexing phase by clustering.

14 Our future work aims to verify our proposed semi-supervised clustering
15 model with larger image databases such as Corel30k, MIRFLICKR, to prove
16 experimentally the convergence of our algorithm, and to look for different
17 strategies for deducing the pairwise constraints or for representing the clus-
18 tering results that could improve the performance of our model in the context
19 of huge image databases.

20 **References**

- 21 Abdi, H., Williams, L. J., 2010. Principal component analysis.
- 22 Basu, S., Banerjee, A., Mooney, R. J., 2002. Semi-supervised clustering by
23 seeding. In: Proceedings of the Nineteenth International Conference on

- 1 Machine Learning. ICML '02. Morgan Kaufmann Publishers Inc., San
2 Francisco, CA, USA, pp. 27–34.
- 3 Basu, S., Bilenko, M., Mooney, R. J., 2004. A probabilistic framework for
4 semi-supervised clustering. In: Proceedings of the tenth ACM SIGKDD
5 international conference on Knowledge discovery and data mining. KDD
6 '04. ACM, New York, NY, USA, pp. 59–68.
- 7 Bayer, R., McCreight, E. M., 1972. Organization and maintenance of large
8 ordered indexes. *Acta Informatica* 1, 173–189.
- 9 Beckmann, N., Kriegel, H.-P., Schneider, R., Seeger, B., 1990. The r^* -tree: an
10 efficient and robust access method for points and rectangles. In: Proceed-
11 ings of the 1990 ACM SIGMOD international conference on Management
12 of data. SIGMOD '90. ACM, New York, NY, USA, pp. 322–331.
- 13 Bentley, J. L., Sep. 1975. Multidimensional binary search trees used for as-
14 sociative searching. *Commun. ACM* 18 (9), 509–517.
- 15 Berchtold, S., Keim, D. A., Kriegel, H.-P., 1996. The x-tree: An index struc-
16 ture for high-dimensional data. In: Proceedings of the 22th International
17 Conference on Very Large Data Bases. VLDB '96. Morgan Kaufmann Pub-
18 lishers Inc., San Francisco, CA, USA, pp. 28–39.
- 19 Dubey, A., Bhattacharya, I., Godbole, S., 2010. A cluster-level semi-
20 supervision model for interactive clustering. In: Proceedings of the 2010
21 European conference on Machine learning and knowledge discovery in
22 databases: Part I. ECML PKDD'10. Springer-Verlag, Berlin, Heidelberg,
23 pp. 409–424.

- 1 Guttman, A., 1984. R-trees: A dynamic index structure for spatial searching.
2 In: International Conference on Management of Data. ACM, pp. 47–57.
- 3 Henrich, A., Six, H. W., Widmayer, P., 1989. The lsd tree: spatial access to
4 multidimensional and non-point objects. In: Proceedings of the 15th in-
5 ternational conference on Very large data bases. VLDB '89. Morgan Kauf-
6 mann Publishers Inc., San Francisco, CA, USA, pp. 45–53.
- 7 Jain, A. K., Murty, M. N., Flynn, P. J., 1999. Data clustering: a review.
8 ACM Comput. Surv. 31 (3), 264–323.
- 9 Katayama, N., Satoh, S., 1997. The sr-tree: an index structure for high-
10 dimensional nearest neighbor queries. In: Proceedings of the 1997 ACM
11 SIGMOD international conference on Management of data. SIGMOD '97.
12 ACM, New York, NY, USA, pp. 369–380.
- 13 Kulis, B., Basu, S., Dhillon, I., Mooney, R., 2005. Semi-supervised graph
14 clustering: a kernel approach. In: In ICML 05: Proceedings of the 22nd
15 international conference on Machine learning. ACM Press, pp. 457–464.
- 16 Lai, H. P., Visani, M., Boucher, A., Ogier, J.-M., 2012a. An experimental
17 comparison of clustering methods for content-based indexing of large image
18 databases. Pattern Analysis and Applications 15, 345–366.
- 19 Lai, H. P., Visani, M., Boucher, A., Ogier, J.-M., 2012b. Unsupervised and
20 semi-supervised clustering for large image database indexing and retrieval.
21 In: IEEE International Conference on Computing and Communication
22 Technologies, Research, Innovation, and Vision for the Future (RIVF). pp.
23 1–6.

- 1 Lance, G. N., Williams, W. T., 1967. A general theory of classificatory sorting
2 strategies ii. clustering systems. *The computer journal* 10 (3), 271–277.
- 3 Likas, A., Vlassis, N., Verbeek, J. J., Feb. 2003. The global k-means clustering
4 algorithm. *Pattern Recognition* 36 (2), 451–461.
- 5 Lowe, D. G., Nov. 2004. Distinctive image features from scale-invariant key-
6 points. *Int. J. Comput. Vision* 60 (2), 91–110.
7 URL <http://dx.doi.org/10.1023/B:VISI.0000029664.99615.94>
- 8 Nievergelt, J., Hinterberger, H., Sevcik, K. C., 1988. Readings in database
9 systems. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, Ch.
10 The grid file: an adaptable, symmetric multikey file structure, pp. 582–598.
- 11 Pearson, K., 1901. On lines and planes of closest fit to systems of points in
12 space. *Philosophical Magazine* 6 (2), 559–572.
- 13 Robinson, J. T., 1981. The k-d-b-tree: a search structure for large multidimensional
14 dynamic indexes. In: *Proceedings of the 1981 ACM SIGMOD international conference on Management of data. SIGMOD '81*. ACM, New
15 York, NY, USA, pp. 10–18.
- 16
- 17 Rosenberg, A., Hirschberg, J., 2007. V-measure: A conditional entropy-based
18 external cluster evaluation measure. In: *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and*
19 *Computational Natural Language Learning (EMNLP-CoNLL)*. pp. 410–
20 420.
21
- 22 Rousseeuw, P., Nov. 1987. Silhouettes: a graphical aid to the interpretation
23 and validation of cluster analysis. *J. Comput. Appl. Math.* 20 (1), 53–65.

- 1 Sellis, T. K., Roussopoulos, N., Faloutsos, C., 1987. The r+-tree: A dy-
2 namic index for multi-dimensional objects. In: Proceedings of the 13th
3 International Conference on Very Large Data Bases. VLDB '87. Morgan
4 Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 507–518.
- 5 Sivic, J., Zisserman, A., 2003. Video google: a text retrieval approach to
6 object matching in videos. In: Computer Vision, 2003. Proceedings. Ninth
7 IEEE International Conference on. pp. 1470–1477.
- 8 van de Sande, K. E. A., Gevers, T., Snoek, C. G. M., 2008. Evaluation of
9 color descriptors for object and scene recognition. In: IEEE Conference on
10 Computer Vision and Pattern Recognition.
- 11 Wagstaff, K., Cardie, C., Rogers, S., Schrödl, S., 2001. Constrained k-means
12 clustering with background knowledge. In: Proceedings of the Eighteenth
13 International Conference on Machine Learning. ICML '01. Morgan Kauf-
14 mann Publishers Inc., San Francisco, CA, USA, pp. 577–584.
- 15 White, D. A., Jain, R., 1996. Similarity indexing with the ss-tree. In: Pro-
16 ceedings of the Twelfth International Conference on Data Engineering.
17 ICDE '96. IEEE Computer Society, Washington, DC, USA, pp. 516–523.
- 18 Xu, R., Wunsch, D., I., 2005. Survey of clustering algorithms. *Neural Net-*
19 *works, IEEE Transactions on* 16 (3), 645–678.
- 20 Zhang, T., Ramakrishnan, R., Livny, M., 1996. Birch: an efficient data clus-
21 tering method for very large databases. In: Proceedings of the 1996 ACM
22 SIGMOD international conference on Management of data. SIGMOD '96.
23 ACM, New York, NY, USA, pp. 103–114.