

Graph Nodes Clustering based on the Commute-Time Kernel

Luh Yen¹, Francois Fouss¹, Christine Decaestecker²,
Pascal Francq³ & Marco Saerens¹

1- Université catholique de Louvain, ISYS, IAG, Place des Doyens 1,
B-1348 Louvain-la-Neuve, Belgium,

{luh.yen, francois.fouss, marco.saerens}@uclouvain.be

2- Université libre de Bruxelles, Institut de Pharmacie, Boulevard du
Triomphe 205/01, 1050 Bruxelles, Belgium, cdecaes@ulb.ac.be

3- Université libre de Bruxelles, STIC, Av. Fr. Roosevelt 50,
1050 Bruxelles, Belgium, pfrancq@ulb.ac.be

October 10, 2006

Abstract

This work presents a kernel method for clustering the nodes of a weighted, undirected, graph. The algorithm is based on a two-step procedure. First, the sigmoid commute-time kernel (\mathbf{K}_{CT}), providing a similarity measure between any couple of nodes by taking the indirect links into account, is computed from the adjacency matrix of the graph. Then, the nodes of the graph are clustered by performing a kernel k-means or fuzzy k-means on this CT kernel matrix. For this purpose, a new, simple, version of the kernel k-means and the kernel fuzzy k-means is introduced. The joint use of the CT kernel matrix and kernel clustering appears to be quite successful. Indeed, this methodology provides good results, outperforming the spherical k-means, on a document clustering problem involving the newsgroups database.

1 Introduction

This work presents a general methodology for clustering the nodes of a weighted, undirected, graph. Graph nodes clustering is an important issue that has been the subject of much recent work; see for instance [5], [6], [8], [14], [19] and [21].

On the other hand, kernel-based algorithms are characterized by two properties: they allow (i) to compute implicitly similarities in a high-dimensional space where the data are more likely to be well-separated and (ii) to compute

similarities between structured objects that cannot be naturally represented by a simple set of features. The latter property will be illustrated in this paper, with the general objective of clustering nodes of a graph, according to some similarity measure between them. It relies on two independent steps:

1. First, define a kernel matrix from the adjacency matrix of the graph, capturing similarities between nodes;
2. Then, use a kernel k-means or fuzzy k-means in order to cluster the nodes of the graph, i.e. the elements of the kernel matrix.

More precisely, given a weighted, undirected, graph, a kernel matrix defining similarities between the nodes is first computed. These similarities take both direct and indirect links into account; they therefore take the indirect paths between the nodes into consideration, and have the nice property of decreasing when the number of paths connecting two nodes increases and when the “length” of any path decreases (the communication is facilitated). In short, two nodes are considered as similar if there are many short paths connecting them. Such kernel matrices have already been used in collaborative recommendation problems with promising results. Indeed, in [10], seven kernel matrices were investigated.

Based on this kernel matrix, nodes are clustered thanks to a kernel clustering. The kernel clustering algorithms proposed in this paper differ from existing ones ([3], [11], [12], [22], [24] and [25]) by the fact that a prototype is explicitly defined for each cluster. Defining prototypes in each cluster is more natural since it allows to mimic the iterative update rules reminiscent from k-means and fuzzy k-means in the sample space, instead of the feature space. In addition to be very similar to the original k-means and fuzzy k-means algorithms, this sample-based method can easily be extended to variable-metric or multi-prototype kernel k-means, in the same way as the original k-means and fuzzy k-means [7]. In addition to this, the resulting algorithm is very simple and natural.

The performances are evaluated on the problem of clustering newsgroups documents, and compared to the popular spherical k-means algorithm, which is especially designed for document clustering [4]. The collection of documents is viewed as a graph and the basic problem is to cluster the documents in order to eventually retrieve the newsgroups. The results indicate that the introduced algorithms perform well in comparison with the the spherical k-means, with significant improvement.

The paper is organized as follows. Section 2 introduces the sigmoid commute-time kernel (\mathbf{K}_{CT}) on a graph that will be used as similarity measure for clustering the nodes. Section 3 derives our version of the kernel k-means and kernel fuzzy k-means, while Section 4 shows the results obtained on the newsgroups database. Section 5 is the conclusion.

2 The sigmoid commute-time kernel on a graph

2.1 Basic notations and definitions

Let us consider that we are given a weighted, undirected, graph, G , with symmetric weights $w_{ij} > 0$ between every couple of nodes, i and j , which are linked by an edge (say G has n nodes in total). The weight w_{ij} of the edge connecting node i and node j should be set to some meaningful value, with the following convention: the more important the relation between node i and node j , the larger the value of w_{ij} , and consequently the easier the communication through the edge. The elements a_{ij} of the adjacency matrix \mathbf{A} of the graph are defined in a standard way as $a_{ij} = w_{ij}$ if node i is connected to node j and 0 otherwise. Based on the adjacency matrix, the Laplacian matrix \mathbf{L} of the graph is defined in the usual manner: $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where $\mathbf{D} = \text{Diag}(a_{i.})$ is the degree matrix, with diagonal entries $d_{ii} = [\mathbf{D}]_{ii} = a_{i.} = \sum_{j=1}^n a_{ij}$. Furthermore, the volume of the graph is defined as $V_G = \text{vol}(G) = \sum_{i=1}^n d_{ii} = \sum_{i,j=1}^n a_{ij}$. We suppose that the graph has a single connected component; that is, any node can be reached from any other node of the graph. In this case, \mathbf{L} has rank $n - 1$, where n is the number of nodes [2]. Moreover, it can be shown that \mathbf{L} is symmetric and positive semidefinite (see for instance [2]).

2.2 The sigmoid commute time kernel

The ‘‘commute time’’ kernel [16], [9] takes its name from the **average commute time**, $n(i, j)$, which is defined as the average number of steps a random walker, starting in node $i \neq j$, will take before entering a node j for the first time, and go back to i . Indeed, we associate a Markov chain to the graph in the following obvious manner. A state is associated to every node (n in total), and the transition probabilities are given by $p_{ij} = a_{ij}/a_{i.}$ where $a_{i.} = \sum_{j=1}^n a_{ij}$. One can show [16], [9] that, in this case, the average commute time can be computed thanks to

$$n(i, j) = V_G (\mathbf{e}_i - \mathbf{e}_j)^T \mathbf{L}^+ (\mathbf{e}_i - \mathbf{e}_j) \quad (1)$$

where every node i of the graph is represented by a basis vector, \mathbf{e}_i (the i -th column of the identity matrix \mathbf{I}), in the Euclidean space \mathfrak{R}^n and V_G is the volume of the graph. \mathbf{L}^+ is the Moore-Penrose pseudoinverse of the Laplacian matrix of the graph and is positive semidefinite. Thus, (1) is a **Mahalanobis distance** between the nodes of the graph and is referred to as the ‘‘commute time distance’’ or the ‘‘resistance distance’’ because of a close analogy with the effective resistance in electrical networks [9].

One can further show that \mathbf{L}^+ is the matrix containing the inner products of the node vectors in the Euclidean space where these node vectors are exactly separated by commute time distances. In other words, the entries of \mathbf{L}^+ can be viewed as similarities between nodes and \mathbf{L}^+ can be considered as a kernel matrix:

$$\mathbf{K} = \mathbf{L}^+ \quad (2)$$

The **sigmoid commute time kernel** \mathbf{K}_{CT} is obtained by applying a sigmoid transform [17] on \mathbf{K} . In other words, each element of the kernel matrix is given by the formula

$$[\mathbf{K}_{CT}]_{ij} = 1/(1 + \exp[a l_{ij}^+/\sigma]) \quad (3)$$

where $l_{ij}^+ = [\mathbf{L}^+]_{ij}$ and σ a normalizing factor, corresponding to the standard deviation of the elements of \mathbf{L}^+ . The parameter a will be set to a constant value determined by informal preliminary tests. The sigmoid function aims to normalize the range of the similarities in the interval $[0, 1]$ [17]. Notice, however, that the resulting matrix is not necessarily positive semi-definite so that, strictly speaking, it is not a kernel matrix.

3 Kernel k-means and fuzzy k-means

We now introduce our kernel version of two well-known clustering algorithms.

3.1 Kernel k-means

The goal is to design an iterative algorithm aiming to minimize a cost function which, in the case of a standard k-means, can be defined, in the feature space, as the total within-cluster inertia:

$$J(\mathbf{g}_1, \dots, \mathbf{g}_m) = \sum_{k=1}^m \sum_{i \in C_k} \|\mathbf{x}_i - \mathbf{g}_k\|^2 \quad (4)$$

where the first sum is taken on the m clusters, while the second sum is taken on the nodes i belonging to cluster k , $i \in C_k$. In Equation (4), \mathbf{x}_i is the feature vector corresponding to node i , \mathbf{g}_k is a prototype vector of cluster k and $\|\mathbf{x}_i - \mathbf{g}_k\|$ is the Euclidean distance between the node vector and the cluster prototype it belongs to. The number of clusters, m , is provided a priori by the user.

We denote by \mathbf{X} the data matrix containing the transposed node vectors as rows, that is, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$. Let us now define the following change of parameter: $\mathbf{g}_k \rightarrow \mathbf{X}^T \boldsymbol{\gamma}_k$, corresponding to the ‘kernel trick’ (see [18]). It aims to express the prototype vectors, \mathbf{g}_k , in terms of the node vectors, \mathbf{x}_i (the columns of \mathbf{X}^T). Now, recompute the within-class inertia in terms of the $\boldsymbol{\gamma}_k$ and the inner products:

$$\begin{aligned} J(\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_m) &= \sum_{k=1}^m \sum_{i \in C_k} (\mathbf{x}_i - \mathbf{g}_k)^T (\mathbf{x}_i - \mathbf{g}_k) \\ &= \sum_{k=1}^m \sum_{i \in C_k} (\mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{g}_k + \mathbf{g}_k^T \mathbf{g}_k) \\ &= \sum_{k=1}^m \sum_{i \in C_k} (k_{ii} - 2\mathbf{k}_i^T \boldsymbol{\gamma}_k + \boldsymbol{\gamma}_k^T \mathbf{K} \boldsymbol{\gamma}_k) \end{aligned} \quad (5)$$

where $\mathbf{K} = \mathbf{X}\mathbf{X}^T$, $k_{ii} = [\mathbf{K}]_{ii} = \mathbf{x}_i^T \mathbf{x}_i$, $\mathbf{k}_i = \mathbf{X}\mathbf{x}_i = \text{col}_i(\mathbf{K})$.

The k-means iteratively minimizes J by proceeding in two steps, (1) re-allocation of the node vectors while keeping the prototype vectors fixed, and (2) re-computation of the prototype vectors, $\boldsymbol{\gamma}_k$, while maintaining the cluster labels of the nodes fixed.

Clearly, the re-allocation step minimizing J is

$$l_i = \arg \min_k \{ \boldsymbol{\gamma}_k^T \mathbf{K} \boldsymbol{\gamma}_k - 2\mathbf{k}_i^T \boldsymbol{\gamma}_k \} \quad (6)$$

where l_i contains the cluster label of node i .

For the computation of the prototype vector, by taking the gradient of J with respect to $\boldsymbol{\gamma}_k$ and setting the result equal to $\mathbf{0}$, we obtain

$$\mathbf{K}\boldsymbol{\gamma}_k = \frac{1}{n_k} \sum_{i \in C_k} \mathbf{k}_i \quad (7)$$

where n_k is the number of nodes belonging to cluster k . By looking carefully to Equation (7), we immediately observe from the left-hand side of (7) that $\mathbf{K}\boldsymbol{\gamma}_k$ is a linear combination of the \mathbf{k}_i , while the right-hand side of (7) is also a linear combination of the \mathbf{k}_i . Therefore, one solution to this linear system of equations is simply the following: $\boldsymbol{\gamma}_k$ contains $1/n_k$ if $i \in C_k$ and 0 otherwise. Therefore the prototype recomputation step is

$$\gamma_{ki} = [\boldsymbol{\gamma}_k]_i = \begin{cases} 1/n_k & \text{if node } i \in C_k \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

This two-step procedure (equations (6) and (8)) is iterated until convergence.

3.2 Kernel fuzzy k-means

We now apply the same procedure for deriving a kernel fuzzy k-means. This time, the cost function is

$$J(\mathbf{g}_1, \dots, \mathbf{g}_m; \mathbf{U}) = \sum_{k=1}^m \sum_{i=1}^n u_{ik} \|\mathbf{x}_i - \mathbf{g}_k\|^2 \text{ with } \sum_{k=1}^m u_{ik}^{1/q} = 1 \text{ for all } i \quad (9)$$

where the u_{ik} define the degree of membership of node i to cluster C_k . The parameter $q > 1$ is controlling the degree of fuzzyness of the membership functions. As for the kernel k-means, we perform the change of parameter $\mathbf{g}_k \rightarrow \mathbf{X}^T \boldsymbol{\gamma}_k$, leading to

$$J(\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_m; \mathbf{U}) = \sum_{k=1}^m \sum_{i=1}^n u_{ik} (k_{ii} - 2\mathbf{k}_i^T \boldsymbol{\gamma}_k + \boldsymbol{\gamma}_k^T \mathbf{K} \boldsymbol{\gamma}_k) \quad (10)$$

We thus introduce the following Lagrange function

$$L(\gamma_1, \dots, \gamma_m; \mathbf{U}) = \sum_{k=1}^m \sum_{i=1}^n u_{ik} (k_{ii} - 2\mathbf{k}_i^T \gamma_k + \gamma_k^T \mathbf{K} \gamma_k) + \sum_{i=1}^n \lambda_i \left[\sum_{k=1}^m u_{ik}^{1/q} - 1 \right] \quad (11)$$

Taking the gradient of L with respect to u_{ik} and setting the result equal to zero allows to compute the membership function,

$$u_{ik} = \left[\frac{(k_{ii} - 2\mathbf{k}_i^T \gamma_k + \gamma_k^T \mathbf{K} \gamma_k)^{-1/(q-1)}}{\sum_{l=1}^m (k_{ii} - 2\mathbf{k}_i^T \gamma_l + \gamma_l^T \mathbf{K} \gamma_l)^{-1/(q-1)}} \right]^q \quad (12)$$

Moreover, taking the gradient with respect to γ_k and setting the result to zero provides

$$\mathbf{K} \gamma_k = \frac{\sum_{i=1}^n u_{ik} \mathbf{k}_i}{\sum_{i=1}^n u_{ik}} \quad (13)$$

Thus, the re-computation of the prototype vectors is simply

$$\gamma_{ki} = [\gamma_k]_i = \frac{u_{ik}}{\sum_{i=1}^n u_{ik}} \quad (14)$$

The equations (12) and (14) are iterated until convergence.

4 Experiments

4.1 Data set

In order to test the performances of the \mathbf{K}_{CT} k-means and the \mathbf{K}_{CT} fuzzy k-means, both algorithms will be assessed on a real data set and compared to the spherical k-means [4]. Our experiments were performed on the newsgroups data set [13]; it is composed of 20,000 unstructured documents, taken from 20 discussion groups (newsgroups) of the Usenet diffusion list. As the data set is composed of documents, the clustering performances of both methods will be compared to those of the spherical k-means [4], which is a reference in text mining.

For our experiment, 9 subsets including different topics are extracted from the original database, as listed in figure 4.1. More precisely, for each subset

200 documents are sampled from different newsgroups. Thus, the three first subsets (G-2cl-A, G-2cl-B, G-2cl-C) contain 400 documents sampled from two newsgroups topics, the next three subsets (G-3cl-A, G-3cl-B, G-3cl-C) contain 600 documents sampled from three topics and the last three subsets (G-5cl-A, G-5cl-B, G-5cl-C) contain 1000 documents sampled from five topics. The selected topics can be well separated or related such as computer/graphics and computer/pchardware in subset G-5cl-B. Both the classification rate (obtained by comparing the clustering to the real newsgroups and performing an optimal assignment) and the adjusted Rand index (with values scaled in $[0, 1]$; see for instance [7]) will be reported.

G-2cl-A	politics/general, sport/baseball
G-2cl-B	computer/graphics, motor/motorcycles
G-2cl-C	space/general, politics/mideast
G-3cl-A	sport/baseball, space/general, politics/mideast
G-3cl-B	computer/windows, motor/autos, religion/general
G-3cl-C	sport/hockey, religion/atheism, medicine/general
G-5cl-A	computer/graphics, motor/motorcycles, politics/mideast, space/general, sport/baseball
G-5cl-B	computer/graphics, computer/pchardware, motor/autos, religion/atheism, politics/mideast
G-5cl-C	computer/machardware, sport/hockey, medicine/general, religion/general, forsale/general

Figure 1: Document subsets used in our experiments. Nine subsets are selected from the *Newsgroups* data set, with 2, 3 or 5 topics. For each subset, 200 documents are randomly selected from each topic.

4.2 Graph definition

The newsgroups data set can be seen as a large bipartite graph between documents and terms. Each document node is connected to terms nodes with each edge weighted by *tf.idf* [20]. After some preprocessing steps (see below) aiming to reduce the number of terms, a graph involving only documents is computed from this bipartite graph in the following way: the link between two documents is given by the sum of all document-term-document paths connecting them and passing through the terms they have in common. In other words, if \mathbf{W} represents the term-document matrix containing the *tf.idf* factors, the adjacency matrix of the resulting document-document graph is provided by $\mathbf{A} = \mathbf{W}^T \mathbf{W}$.

4.3 Preprocessing steps

In order to reduce the high dimensionality of the feature space (terms), the following standard preprocessing steps are performed on the data set before the clustering experiment.

1. Stopwords without useful information are eliminated.
2. Porter’s stemming algorithm [15] is applied so that each word is reduced to its ‘root’.
3. Words that occur too few times (< 3) or in too few documents (< 2) are considered as no content-bearing and are eliminated.
4. Mutual information between terms and documents is computed. For a word y , the mutual information with the documents of the data set [23] is given by

$$I(y) = \sum_x \log p(x, y) / p(x)p(y), \quad (15)$$

where x is the document of the data set. Words with a small value of mutual information (fixed at 20% of $I(y)$ ’s median) are eliminated.

5. The term-document matrix \mathbf{W} is constructed with the remaining words and documents. Element $[\mathbf{W}]_{ij}$ of the matrix contains the value of *tf.idf* factor between the term i and the document j .
6. Each row of the term-document matrix \mathbf{W} is normalized to 1.

For example, the subset G-2cl-A is composed of 400 documents, and 2898 terms with stopwords already eliminated. After preprocessing, only 1490 terms are kept. Finally, the adjacency matrix of the documents graph \mathbf{A} is given by the document-document matrix $\mathbf{W}^T\mathbf{W}$. Based on \mathbf{A} , \mathbf{K}_{CT} is computed by Equation (3).

4.4 Experimental settings

Suppose we have a graph of n nodes to be partitioned into m clusters. First, the prototype vectors γ_i ($i = 1, \dots, m$) are initialized by randomly selecting m columns of the identity matrix \mathbf{I} . Then, each algorithm is launched 30 times (30 runs), and the classification rate as well as the adjusted Rand index, averaged on the 30 runs, are computed. \mathbf{K}_{CT} k-means, \mathbf{K}_{CT} fuzzy k-means are run with the document-document matrix \mathbf{A} , fuzzy k-means and spherical k-means is run with the term-document matrix \mathbf{W} after preprocessing.

Each run of a clustering algorithm consists in 50 trials: the clustering algorithm is launched 50 times and the best solution among the 50 trials, having the minimal within-class inertia, is sent back as the solution.

Two parameters need to be tuned. The first one is the parameter a for computing the sigmoid transform of the \mathbf{K}_{CT} (see Equation (3)). The second

one is the parameter q which controls the degree of fuzzyness for the kernel fuzzy k-means (see Equation (12)). Based on preliminary informal experiment, the parameters a and q were set to 7 and 1.2 respectively, for all experiments.

4.5 Experimental results and discussion

The results (the classification rate as well as the adjusted Rand index, each averaged on 30 runs) of the three clustering algorithms (\mathbf{K}_{CT} k-means, \mathbf{K}_{CT} fuzzy k-means and spherical k-means) on the nine document subsets are reported in Table 1.

	\mathbf{K}_{CT} k-means		\mathbf{K}_{CT} fuzzy k-means		Sph. k-means	
	class. rate	adj. Rand	class. rate	adj. Rand	class. rate	adj. Rand
G-2cl-A	97.25 %	0.95	97.25 %	0.95	91.76 %	0.85
G-2cl-B	91.23 %	0.84	91.46 %	0.84	81.46 %	0.70
G-2cl-C	95.71 %	0.92	95.99 %	0.92	94.82 %	0.90
G-3cl-A	94.42 %	0.92	94.83 %	0.93	89.23 %	0.85
G-3cl-B	93.37 %	0.91	93.14 %	0.90	86.71 %	0.82
G-3cl-C	93.65 %	0.91	92.44 %	0.89	87.35 %	0.83
G-5cl-A	82.64 %	0.80	86.49 %	0.84	75.29 %	0.74
G-5cl-B	70.75 %	0.75	77.08 %	0.78	64.43 %	0.69
G-5cl-C	77.12 %	0.76	79.54 %	0.78	65.23 %	0.69

Table 1: Comparison of the clustering performances (classification rate in % and adjusted Rand index with value scaled in $[0, 1]$) for the \mathbf{K}_{CT} k-means, \mathbf{K}_{CT} fuzzy k-means and spherical k-means.

We observe that the \mathbf{K}_{CT} k-means and the \mathbf{K}_{CT} fuzzy k-means outperform the spherical k-means on the nine subsets. Moreover, the \mathbf{K}_{CT} fuzzy k-means provides slightly better results than the two other methods. This can be partly explained by the fact that the newsgroups data set is fuzzy itself, as discussed in [1]. It is hard to define clear boundaries between the different topics: a discussion within a specific newsgroup can also be related to other domains. The Figure 2 shows an example of fuzzy clustering for the subset G-5cl-A. For each topic, the membership probability of each document of the subset is displayed. Most misclassified documents belong to the topic space/general and are wrongly assigned to the computer/graphics topic.

A close examination of the membership matrix on Figure 2 shows that, for some documents, the membership value is almost equal for the five clusters obtained on the subset G-5cl-A. An example of such a document is shown here-below. It is document 103155, picked from the topic motor/motorcycles.

Document 103155 : *Date: Mon, 5 Apr 1993 14:49:08 EDT*
From: <LRR105@psuvm.psu.edu>
Message-ID: <93095.144908LRR105@psuvm.psu.edu>

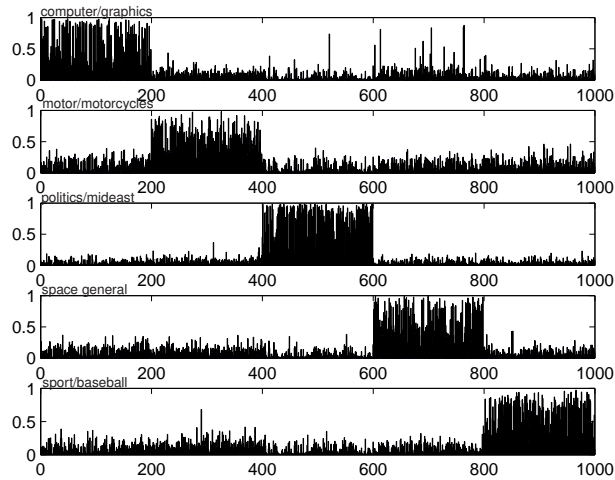


Figure 2: Membership value for each document and each cluster of the G-5cl-A subset.

Newsgroups: rec.autos.tech,rec.motorcycles
Subject: Re: Tools Tools Tools
WHAT IS THE FLANK DRIVE EVERYONES TALKING ABOUT?

The fuzzy membership vector for this document, provided by \mathbf{K}_{CT} fuzzy k-means, indicate that the document belongs to the topic politics/mideast with a membership value of 0.28, space/general with 0.26, computer/graphics with 0.21, motor/motorcycles with 0.13 and sport/baseball with 0.11. Actually, it indicates that the algorithm fails to determine what is the correct topic of this message. Actually, even a person not used to the domain will have difficulty to determine the subject of the discussion.

Consequently, in order to examine the influence of the fuzzy membership function, we decided to set the label of a document to ‘undefined’ when the maximal membership value for this document is below a given certainty threshold. Figure 3(a) shows the classification rate in terms of this certainty threshold, and Figure 3(b) shows the number of unclassified documents, both for the subsets G-5cl-A (upper figures) and G-5cl-C (lower figures). With the certainty threshold fixed at 0.45, the classification rate of G-5cl-A is improved to 98.2% but, in this case, 32% of the documents are unclassified.

5 Conclusions and further work

We introduced a new method allowing to cluster the nodes of a weighted graph by exploiting the links between them. It is based on a recently introduced kernel on a graph, the commute-time kernel, combined with a kernel clustering. The

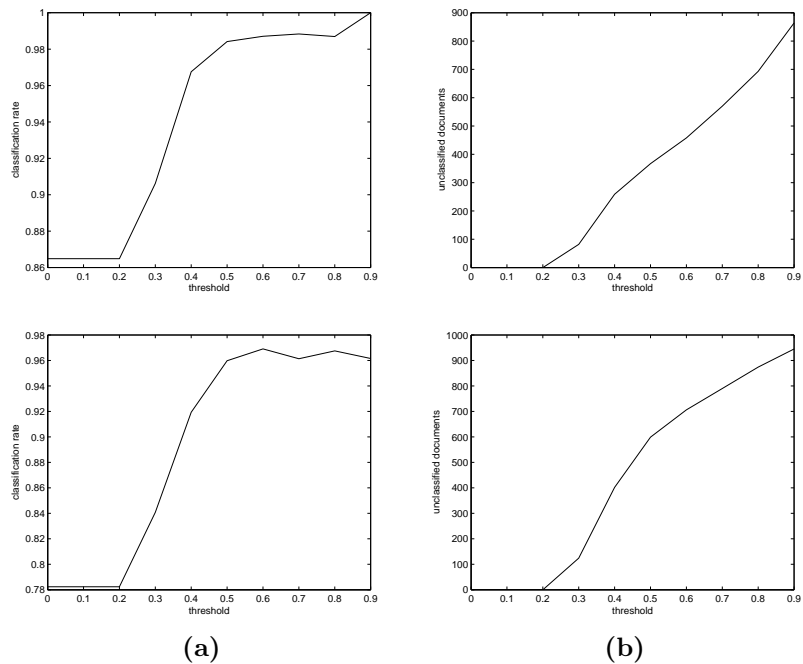


Figure 3: \mathbf{K}_{CT} fuzzy clustering results (classification rate) in terms of the membership certainty threshold for subsets G-5cl-A and G-5cl-C: (a) Classification rate evolution in terms of the membership threshold value; (b) Number of unclassified documents in terms of the membership threshold.

obtained results are promising since the proposed methodology outperforms the standard spherical k-means on a difficult document clustering problem. Further work will be devoted to (1) additional experiments on biological as well as text databases, and (2) developing kernel versions of the Gaussian mixture and the entropy-based fuzzy clustering and assessing their performances.

References

- [1] M. W. Berry, editor. *Survey of Text Mining*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2003.
- [2] F. R. Chung. *Spectral graph theory*. American Mathematical Society, 1997.
- [3] I. S. Dhillon, Y. Guan, and B. Kulis. Kernel k-means, spectral clustering and normalized cuts. In *Proceedings of the 2004 ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 551–556. ACM Press, 2004.
- [4] I. S. Dhillon and D. S. Modha. Concept decompositions for large sparse text data using clustering. *Machine Learning*, 42(1):143–175, 2001.
- [5] C. Ding and X. He. Linearized cluster assignment via spectral ordering. In *ICML '04: Proceedings of the twenty-first international conference on Machine learning*, page 30, New York, NY, USA, 2004. ACM Press.

- [6] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. *Machine Learning*, 56(1-3):9–33, 2004.
- [7] B. S. Everitt, S. Landau, and M. Leese. *Cluster Analysis*. Arnold Publishers, 2001.
- [8] G. W. Flake, R. E. Tarjan, and K. Tsioutsoulis. Graph clustering and minimum cut trees. *Internet Math*, 1(4):385–408, 2003.
- [9] F. Fouss, A. Pirotte, J.-M. Renders, and M. Saerens. Random-walk computation of similarities between nodes of a graph, with application to collaborative recommendation. *Submitted to the IEEE Transactions on Knowledge and Data Engineering*, 2005.
- [10] F. Fouss, L. Yen, A. Pirotte, and M. Saerens. An experimental investigation of graph kernels on two collaborative recommendation tasks. To appear in the proceedings of the 2006 IEEE International Conference on Data Mining (ICDM 2006).
- [11] M. Girolami. Mercer kernel-based clustering in feature space. *IEEE Transactions on Neural Networks*, 13(3):780–784, May 2002.
- [12] D.-W. Kim, K. Y. Lee, D. Lee, and K. H. Lee. Evaluation of the performance of clustering algorithms in kernel-induced feature space. *Pattern Recognition*, 38(4):607–611, 2005.
- [13] K. Lang. 20 newsgroup dataset. Available from <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/news20.html>.
- [14] M. Newman. Detecting community structure in networks. *The European Physical Journal B*, 38:321–330, 2004.
- [15] M. F. Porter. An algorithm for suffix stripping. *Program*, 14(3):130–137, 1980.
- [16] M. Saerens, F. Fouss, L. Yen, and P. Dupont. The principal components analysis of a graph, and its relationships to spectral clustering. *Proceedings of the 15th European Conference on Machine Learning (ECML 2004). Lecture Notes in Artificial Intelligence, Vol. 3201, Springer-Verlag, Berlin*, pages 371–383, 2004.
- [17] B. Scholkopf and A. Smola. *Learning with kernels*. The MIT Press, 2002.
- [18] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [19] S. van Dongen. *Graph Clustering by Flow Simulation*. PhD thesis, University of Utrecht, 2000.
- [20] S. Weiss, N. Indurkha, T. Zhang, and F. Damerau. *Text Mining: Predictive Methods for Analyzing Unstructured Information*. Springer, 2004.
- [21] S. White and P. Smyth. A spectral clustering approach to finding communities in graph. In *SDM*, 2005.
- [22] Z.-D. Wu, W.-X. Xie, and J.-P. Yu. Fuzzy c-means clustering algorithm based on kernel method. In *ICCIMA '03: Proceedings of the 5th International Conference on Computational Intelligence and Multimedia Applications*, page 49, Washington, DC, USA, 2003. IEEE Computer Society.
- [23] H. Zha, X. He, C. H. Q. Ding, M. Gu, and H. D. Simon. Bipartite graph partitioning and data clustering. In *Proc. of ACM 10th Int'l Conf. Information and Knowledge Management (CIKM 2001)*, pages 25–32, 2001.
- [24] D.-Q. Zhang and S.-C. Chen. Fuzzy clustering using kernel method. In *Proceedings of the 2002 International Conference on Control and Automation, 2002. ICCA*, pages 162–163, 2002.
- [25] D.-Q. Zhang and S.-C. Chen. A novel kernelized fuzzy c-means algorithm with application in medical image segmentation. *Artificial Intelligence in Medicine*, 32(1):37–50, 2004.