

DOCUMENT MANAGEMENT USING CLUSTERING ALGORITHMS

CREȚULESCU Radu George¹, PITIC Antoniu Gabriel², MORARIU Daniel³

Lucian Blaga University of Sibiu, Romania

Abstract Document management systems are complex systems, which offer services as storage, versioning, metadata, security, as well as indexing and retrieval capabilities. Large numbers of documents could be automatically grouped into classes of documents, which contain similar information. Therefor we propose to use clustering methods in order to group the documents.

Clustering is an important process in text mining used for groping documents based on their contents in order to extract knowledge. In this paper we will present some requirements for clustering algorithms for a document management system

Keywords: Management, Document Management, Clustering, Cluster Validation

JEL classification: Y80

1. Introduction

In the recent years, significant increases in using documents from the Web and the improvements of the quality and speed of the Internet have transformed our society into one that depends strongly on the information. The huge amount of data that is generated by the process of intercommunication represents important information that accumulates daily and that is stored in form of text documents, databases etc.

¹ lecturer Ph.D., Faculty of Engineering, Computer Sciences Department, Lucian Blaga University of Sibiu, Romania, radu.kretzulescu@ulbsibiu.ro

² teaching assistant Ph.D.c., Faculty of Engineering, Computer Sciences Department, Lucian Blaga University of Sibiu, Romania, antoniu.pitic@ulbsibiu.ro

³ associate professor, Ph.D., Faculty of Engineering, Computer Sciences Department, Lucian Blaga University of Sibiu, Romania, daniel.morariu@ulbsibiu.ro

The document management system (DMS) becomes very important for managing documents.

The elements of a DMS for documents are: metadata extraction integration, capture, validation, indexing, storage, retrieval, distribution, security, workflow, collaboration, versioning, searching, publishing, and reproduction. Also some new systems provide text retrieving and summarization facilities

As mentioned in Berkhin (2006), Han (2001), machine learning software provides the basic techniques for data mining by extracting information from raw data contained in databases. The process usually goes through the following steps:

- transforming the data into a suitable format
- data cleaning
- deduction or conclusions on the extracted data.

Machine learning techniques are divided into two sub domains: supervised learning and unsupervised learning. Under the category of unsupervised learning, one of the main tools is data clustering. This paper attempts to provide taxonomy of the most important algorithms used for clustering. For each algorithm category, we have selected the most common version of the entire family. Below we will present algorithms used in context of document clustering

2. Unsupervised versus supervised learning

In supervised learning, the algorithm receives data (the text documents) and the class label for the corresponding classes of the documents (called labeled data). The purpose of supervised learning is to learn the concepts that correctly classify documents for given classification algorithm. Based on this learning the classifier will be able to predict the correct class for unseen examples. Under this paradigm, it is also possible the appearance of the over-fitting effects. This will happen when the algorithm memorizes all the labels for each case.

The outcomes of supervised learning are usually assessed on a disjoint test set of examples from the training set examples. Classification methods used are varied, ranging from traditional statistical approaches, neural networks to kernel type algorithms Burges (1998).

The quality measure for classification is given by the accuracy of

classification.

In unsupervised learning the algorithm receives only data without the class label (called unlabeled data) and the algorithm task is to find an adequate representation of data distribution.

The central point of this paper is to present the clustering as a key aspect in unsupervised learning.

Some researchers have combined unsupervised and supervised learning that has emerged the concept of semi-supervised learning Benett (1998). In this approach is applied initially an unknown data set in order to make some assumptions about data distribution and then this hypothesis is confirmed or rejected by a supervised approach.

3. Cluster Analysis

Cluster analysis is an iterative process of clustering and cluster validation facilitated by clustering algorithms and cluster validation methods. Cluster analysis includes two major aspects: clustering and cluster validation.

Clustering refers to grouping objects according to certain criteria. To achieve this goal researchers have developed many algorithms Berkin (2006), Jain (1988), Jain (1999), Kaufman (1990). Since there are no general algorithms that can be applied to all types of cases it becomes necessary to apply a validation mechanism so that the user will find an algorithm suitable for its particular case.

Cluster analysis is a process of discovery through exploration. It can be used to discover structures without the need for interpretation Jain (1988).

The validation of the clusters becomes a cluster quality evaluation process.

3.1 Definitions

Def. 3.1: Clustering is the process of grouping physical or abstract objects into classes of similar properties. Jain (1999)

A cluster is a collection of objects that are similar to each other from the same group and are dissimilar to objects that are from different groups.

Def. 3.2 Conceptual clustering – is a process of grouping where the objects in a group will form a class only if they can be described by one concept. This approach differs from conventional clustering where the dissimilarity measure is based on mathematical distances.

The basic idea for clustering, which is the starting point for both previous

presented approaches, is to find those clusters with high inter-cluster similarity and very small intra-cluster similarity.

The main directions of clustering research are focused more on distance-based cluster analysis. Several algorithms have been developed such as partitional algorithms like k-Means and k-Medoids, or hierarchical algorithms. Further investigations have demonstrated the usefulness of other approaches such as algorithms based on suffix trees Meyer (2005), Zamir (1998), the bio-inspired algorithms, such as those based on the behavior of ants in search of food and creation of cemeteries Dorigo (2000), Labroche (2003), particles swarms (particle swarm optimization) Merve (2003) and ontologies Hoto (2003).

3.2. General requirements for clustering algorithms

Typical requirements for clustering algorithms in data mining and text documents Han (2001) are:

- a. Scalability of the algorithms: Many clustering algorithms work very well with small data sets. However, large databases containing millions of objects and using a large sample of these sets would lead to inconclusive results.
- b. Ability to use different types of attributes: Many clustering algorithms use numerical data as input. In some cases it is necessary to apply clustering algorithms on string data, binary data, ordinal data or a combination thereof.
- c. Clusters of arbitrary shape recognition: Many algorithms find clusters using geometric distances such as Euclidean distance or Manhattan distance. These algorithms, however, tend to find spherical clusters with similar size and density. In reality, clusters can be in any form.
- d. Setting input parameters based on minimal knowledge of the field: Many clustering algorithms require certain parameters such as the number of clusters that need to be determined (e.g. *k*-means). The result of clustering can be significantly different depending on the input parameters set. The input parameters for data sets containing high dimensional objects are difficult to be determined.
- e. Ability to use data that contain noise: Many real data sets containing missing data, unknown or outlier. Some algorithms are sensitive to such data and clustering quality becomes poor.

- f. Insensitivity to the order of data processing: Some clustering algorithms may produce different results for different order of the inputs (e.g., Single Pass, BIRCH- Balanced Iterative Reducing and Clustering).
- g. High dimensionality: Data may contain a lot of dimensions and attributes. Many clustering algorithms fail to produce good results even for data with small dimensions. The human eye fails to assess the quality of a cluster by up to three dimensions.
- h. Interpretability and usability: Users want that the clustering results shall be useful, interpretable and understandable.
- i.

4. Data used in clustering

4.1 Data representation

A data matrix represents for the clustering algorithm a matrix with n objects and each has p attributes. The representation of the data will be a matrix of $n \times p$ attributes.

$$\begin{bmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{ij} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix} \quad (1)$$

4.2 Dissimilarity Matrix

This matrix will be a $N \times N$ square matrix and contains the dissimilarity measures of all pairs of objects for the clustering. Since $dis(i, j) = dis(j, i)$ and $dis(i, i) = 0$ we obtain the following dissimilarity matrix

$$\begin{bmatrix} 0 & & & & \\ dis(2,1) & 0 & & & \\ dis(3,1) & dis(3,2) & 0 & & \\ \dots & \dots & \dots & 0 & \\ dis(n,1) & & & dis(n,n-1) & 0 \end{bmatrix} \quad (2)$$

where $dis(i, j)$ is the measure of dissimilarity between two objects.

Objects are clustered according to the similarity between them or by their dissimilarity.

4.3 Dissimilarity versus similarity

For simplifying the notation we not the dissimilarity with $d(i, j)$ is a which is positive number close to 0 when the documents i and j are close to each other and increases when the distance between i and j increases. The dissimilarity between two objects can be obtained by subjective evaluation made by experts based on direct observation or can be calculated based on correlation coefficients.

The similarity $s(i, j)$ is a positive number close to 0 when the two documents are not similar and increases when the documents are more similar.

If the similarity and the dissimilarity coefficients are in the range $[0, 1]$ the following equation is true:

$$d(i, j) = 1 - s(i, j) \tag{3}$$

4.4. Common formula used for dissimilarity/similarity

To calculate the dissimilarity/similarity of objects we calculate the distance between every two objects. The distances must satisfy the following properties:

$d(i, j) \geq 0$	Non-negativity
$d(i, i) = 0$	Identity
$d(i, j) = d(j, i)$	Symmetry
$d(i, j) \leq d(i, h) + d(h, j)$	Triangle rule

Commonly used formulas for dissimilarity/distance are:

a. *Euclidean distance*:

$$d_E(\bar{x}_i, \bar{x}_j) = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2} \tag{3}$$

b. *Manhattan distance*

$$d_{Ma}(\bar{x}_i, \bar{x}_j) = \sum_{k=1}^p |x_{ik} - x_{jk}| \tag{4}$$

c. Minkowski distance

$$d_{Mi}(\bar{x}_i, \bar{x}_j) = \sqrt[q]{\sum_{k=1}^p (x_{ik} - x_{jk})^q} \quad (5)$$

d. Cosine distance

$$d_{\cos}(\bar{x}_i, \bar{x}_j) = \frac{\sqrt{\sum_{k=1}^p x_{ik}^2} \cdot \sqrt{\sum_{k=1}^p x_{jk}^2}}{\sum_{k=1}^p (x_{ik} \cdot x_{jk})} \quad (6)$$

e. Canberra distance

$$d_{CAN}(\bar{x}_i, \bar{x}_j) = \sum_{k=1}^n \frac{|x_{ik} - x_{jk}|}{|x_{ik}| + |x_{jk}|} \quad (7)$$

5. Some clustering algorithms

Arranging data in different groups can be made using different strategies. Based on this, the clustering strategy can be divided into several categories: techniques based on data partitioning, clustering techniques based on hierarchical methods Berkin (2006), methods based on attributes order, methods based on density, grid-based methods and methods based on models. Han (2001).

5.1 Partitioning Methods

If n objects must be grouped into k groups, then a partitioning method constructs k partitions of the objects, each partition is represented by a cluster with $k \leq n$. The clusters are formed taking into account the optimization of a criterion function. This function expresses the dissimilarity between the objects, so that the objects that are grouped into a cluster are similar and objects from different clusters are dissimilar. For this type of grouping method the clusters must satisfy two conditions:

- Each cluster must contain at least one object;

- Each object must be included in a single cluster.

The basic idea of this type of methods is that the algorithm initially starts with a given number k groups representing the number of partitions (clusters) and then it applies a partitioning method that recalculates the k clusters. The method also uses an iterative relocation technique that attempts to improve the partitioning by moving objects from one group to another. The moving criterion needs to respect the condition that the objects of the same cluster are similar (close). There are different criteria to evaluate the quality of clusters. These will be presented in Section 6.

We can identify three different types of partitioning algorithms:

- K-Means clustering algorithms; for this type of algorithm each cluster is represented by the average value of all objects from the same cluster.
- K-Medoids clustering algorithms: for this type of algorithm each cluster is represented by the objects which are closest to the center (medoid) of the cluster.
- Probabilistic algorithms for clustering: a probabilistic method assumes that the data come from a mixture of populations whose distributions and probabilities must be determined.

These algorithms can be used in collections of small to medium data when the resulted clusters can be found in spherical forms.

5.1.1.K-Means algorithm

K-Means algorithm Hartigan (1997), Hartigan (1975), MacQueen (1967) is the most popular algorithm used in scientific and industrial applications. The name of the algorithm comes from the representation of k -clusters C_j whose weights (means) c_j are calculated as the centroid of the points which are grouped in cluster C_j . The value for the parameter k is set at the start of the algorithm and represent the number of cluster that want to be obtained. The similarity between the items from the same cluster C_j is calculated in relation to the centroid. The error criterion used is defined in the formula

$$E(C) = \sum_{j=1}^k \sum_{x_i \in C_j} \|x_i - c_j\|^2 \quad (8)$$

Where x is the given entry object and c_j the centroid of C_j .

The scope is to achieve a greater intra-cluster similarity and very small inter-

cluster similarity so that we get k clusters as compact and simultaneously separated as possible. It is obvious that this type of algorithm works with numerical variables. The algorithm performs the following steps:

1. Generate k random centers in n -dimensional space, which represents the initial centroids.
2. For each object compute the distance between it and all the centroids, then the object is assigned to the closest centroid for computing the distance there are several formulas presented in section 4.4.
3. Once all objects have been properly assigned to the centroids, the positions of the k centroids are recalculated as the center of all samples assigned to each centroid individually.
4. Repeat steps 2 and 3 until the centroids no longer change their positions.
5. The objects assigned to the centroids represent the contents of the final clusters.
- 6.

5.1.2 K-Medoids method

K-medoids algorithm is an adaptation of k-Means algorithm. Instead of computing the centroid of each group, the algorithm chooses a representative element called medoid for each cluster at each iteration. The medoid for each group is calculated by finding an item in the cluster that minimizes the sum:

$$\sum_{j \in C_i} d(i, j) \quad (9)$$

Where C_i is the cluster that contains the object i and $d(i, j)$ the distance between object i and object j .

There are two advantages by using existing objects as cluster centers. First, a medoid can describe usefully a cluster. Second, it is not necessary to calculate the distance between the objects for each iteration of the algorithm. The distances can be computed once and then saved in a so called distances matrix.

The steps of the k-Medoids algorithm can be summarized as follows:

1. Choose k objects random to be the original medoids of the clusters.
2. Assign each object to the nearest cluster associated to the medoid.

3. Recompute the positions for the k-medoids.
4. Repeat steps 2 and 3 until the medoids doesn't change.

Kaufman and Rousseeuw presented in Kaufman (1987) the PAM (Partition Around Medoids) algorithm which is an iteratively implementation of the k-Medoids algorithm.

The K-Medoids algorithm is more robust in terms of noise comparing to the k-Means algorithm. Since k-Medoids is working directly with medoids it is not so much influenced by the outliers like the centroid calculation influences the k-Means algorithm. However in terms of computational cost the k-medoids algorithm has a higher cost because for a single iteration the computational cost is $O(k(n-k)^2)$. The K-medoids algorithm is efficient on small data sets. On large data sets, due to its high computational cost, the algorithm becomes quite slow.

5.2 Hierarchical methods

These types of methods provide a hierarchical structure of the set of objects. There are two approaches of hierarchical methods:

Agglomerative -These methods have a "bottom-up" approach. At the beginning of the algorithm each object is a cluster, and then in the next steps the clusters are merged together based on similarity measures creating a hierarchical structure until all clusters are joined into a single cluster or until another stopping condition (given number of clusters, time, etc.) is reached.

Divisive - These methods have a "top-down" approach. First all objects are considered to be contained in a single cluster, then after successive iterations each cluster is divided into smaller clusters until each object is a cluster or until a stopping condition is reached.

A major problem of hierarchical clustering algorithms is that once the division or merging step has been made it cannot be canceled. This issue also represents a major advantage of this type of calculations due to reduction algorithms avoiding calculating the various combinations of possibilities.

5.3 Method based on the order of words

Suffix Tree Clustering (STC)

The algorithm presented in Zamir (1998) does not require a vector representation of objects (documents) like the algorithms presented in the previous sections. The STC algorithm uses the Suffix Tree Document Model (STDM) to represent the documents. This algorithm will take into account the words order, and creates clusters based on words or groups of consecutive words from the documents. Thus can be a major advantage because it takes account of word order in sentences.

A suffix tree of a document d is a compact tree containing all suffixes s of that document. In our case a suffix is a string consisting of one or more words.

Rules for building the suffix tree:

1. Each internal node other than root must have at least two children and each edge leaving a node is labeled with a nonempty substring n .
2. Any two edges that start from the same node cannot start with the same word.

6. Cluster validation

Each clustering algorithm applied on the same set of data will group the data in different ways depending on the similarity metric used. This makes analysis for the efficiency of clustering algorithms very difficult.

To evaluate the performance or quality of clustering algorithms, objective measures must be established. There are three types of quality measures:

- a. external, when there is an a priori knowledge about the clusters (we have pre-labeled data);
- b. internal, which have no information about clusters;
- c. relative assessing group differences between different solutions.

External measures are applied to both classification and clustering algorithms, while internal and relative measures are applied only to the clustering algorithms.

EXTERNAL VALIDATION MEASURES

External validation measures require pre-labeled data sets for clustering analysis. Because clustering data can be done from many points of view and

the labels may vary from these points of view, the comparison is made to the group containing the most data in a specific category.

Some of the external evaluation measures:

Precision is the percentage of retrieved documents which are really relevant to that category. The value for the precision is in the interval $[0,1]$, with 1 as best. For a cluster C_i and a known class S_j we compute the precision:

$$precision(C_i, S_j) = \frac{|C_i \cap S_j|}{|C_i|} \quad (10)$$

Recall is the percentage of documents that are relevant to that category and are indeed grouped into that category. The value for the recall is in the interval $[0,1]$, with 1 as best.

$$recall(C_i, S_j) = \frac{|C_i \cap S_j|}{|S_j|} \quad (11)$$

Accuracy is the percentage of documents that are correctly grouped into categories depending on the labels of documents (needs labeled documents).

Fmeasure is a measure that combines precision and recall site and is calculated according to formula (13)

$$Fmeasure(C_i, S_j) = \frac{2 \cdot precision(C_i, S_j) \cdot recall(C_i, S_j)}{precision(C_i, S_j) + recall(C_i, S_j)} \quad (12)$$

For each class it will select only the cluster with the highest Fmeasure. In the final Fmeasure for overall measure of a clustering solution is weighted by the size of each cluster.

INTERNAL VALIDATION MEASURES

In WWW (2014) are presented four such metrics:

Compactness—This measure expresses how similar are the data from a given cluster

$$compactness(C) = \sum_{i=1}^{n_c} (\bar{c} - \bar{x}_i)^2 \quad (13)$$

Where C is the current cluster, n_c the number of elements in the cluster, c

the cluster centroid and x_i is an element of the cluster. In other words, these metric measures how "close" the documents within a cluster are. The value is in the interval $[0, \infty)$ and as the lower the better is the measure

Separability–This measure between the clusters expresses how dissimilar the clusters are.

$$separability = \sum_{i,j=1}^c \min(dist(\bar{c}_i, \bar{c}_j)) \quad (14)$$

Thus for each cluster, the nearest cluster is determined. Based on this metric we seek clusters that maximize this value, so that the clusters are very dissimilar.

Balance - The balance is computed as the clusters are formed and expresses how “well-balanced” the formed clusters are.

$$balance = \frac{n / k}{\max_{i=1,k}(n_i)} \quad (15)$$

Where n is the total number of documents, n_i the number of documents in the cluster I (the formula takes into consideration the biggest formed cluster), k is the number of formed clusters.

This measure can take values in the interval $[0,1]$, the maximum value 1 is achieved when all clusters have the same number of documents. A value close to 0 is obtained when the number of documents contained in clusters varies greatly.

7. Conclusions

Clustering in text documents is an unsupervised learning method of classifying documents with a certain degree of similarity using different metrics based on distance. The clustering algorithm must identify (must find) the clusters in which documents are grouped and/or some patterns (rules) that separate one group from another group. There is no predefined taxonomy. It is established when the clustering algorithm run to the set of documents.

Cluster evaluation is a next important task. The formulas presented for internal and external evaluation of clusters can give a measure for evaluate the

quality of clustering process. Still it is very hard to compare different clustering results.

8. Acknowledgment

This work was supported by the strategic grant POSDRU/159/1.5/S/133255, Project ID 133255 (2014), co-financed by the European Social Fund within the Sectorial Operational Program Human Resources Development 2007-2013.

9. References

- Berkhin, P., (2006) - "A Survey of Clustering Data Mining Techniques", Kogan, Jacob;Nicholas, Charles; Teboulle, Marc (Eds.) Grouping Multidimensional Data, Springer Press, pp. 25-72
- Burges, C. J. C., (1998) "A tutorial on support vector machines for pattern recognition." *Data Mining and Knowledge Discovery*, 2(2):121-167
- Dorigo, M., Bonabeau, E., Theraulaz, G., (2000) "Ant Algorithms and stigmergy" *Future Generation Computer Systems* Vol.16, 2000.
- Han, J., Kamber, M., (2001) "Data Mining: Concepts and Techniques", Morgan Kaufmann Publishers
- Hartigan, J. A., Wong, M., (1997) Algorithm as 136: A k-means clustering algorithm, *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, London
- Hartigan, J. A., (1975) "Clustering Algorithms", New York: John Wiley & Sons, Inc, 1975
- Jain, A., K., Dubes, R.,C. (1988) "Algorithms for Clustering Data", Prentice Hall, Englewood Cliffs, NJ
- Jain, A. Murty, M. N., Flynn, P. J., (1999) "Data Clustering: A Review", *ACM Computing Surveys*, Vol. 31(3), pp. 264-323
- Kaufman, L. and Rousseeuw, P.J. (1990) "Finding Groups in Data: An Introduction to Cluster Analysis", Wiley-Interscience, New York (Series in Applied Probability and Statistics)

- Kaufman, L., Rousseeuw, P. J., (1987) "Clustering by means of medoids, in Statistical Data Analysis based on the L, Norm", edited by Y. Dodge, Elsevier/North-Holland, Amsterdam
- Labroche, N., Monmarché, N., Venturini G., (2003) "AntClust: Ant Clustering and Web Usage Mining, In Genetic and Evolutionary Computation", GECCO 2003 Lecture Notes in Computer Science, Volume 2723/2003, 201.
- MacQueen, J. B., (1967) "Some Methods for classification and Analysis of Multivariate Observations", Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, University of California Press, 1:281-297
- Meyer, S., Stein, B., Potthast, M., (2005) "The Suffix Tree Document Model Revisited", Proceedings of the I-KNOW 05, 5th International Conference on Knowledge Management, Journal of Universal Computer Science, pp.596-603, Graz
- Zamir, O, Etzoni, O., (1998) "Web Document Clustering: A Feasibility Demonstration", Proceedings of the 21st International ACM SIGIR Conference on Research and Development in Information Retrieval, Melbourne, Australia
- WWW (2014) Metrics for evaluating clustering algorithms-
<http://www.scribd.com/Clustering/d/28924807-> accessed 25.08.2014