

Using the Self-Organizing Map (SOM) Algorithm, as a Prototype E-Content Retrieval Tool

Athanasios S. Drigas and John Vrettaros

Department of Applied Technologies, NCSR "DEMOKRITOS"
Ag. Paraskeui, Greece
{dr, jvr}.imm.demokritos.gr

Abstract. SOM O.D.I.S.S.E.A.S is an intelligent searching tool using the Self-Organizing Map (SOM) algorithm, as a prototype e-content retrieval tool. The proposed searching tool has the ability to adjust and scale into any e-learning system that requires concept-based queries. In the proposed methodology, maps are used for the automatic replacement of the unstructured, the half structured and the multidimensional data of text, in a way that similar entries in the map are represented near between them. The performance and the functionality of the document organization, and the retrieval tool employing the SOM architecture, are also presented. Furthermore, experiments were performed to test the time performance of a learning algorithm used for the direct creation of teams of terms and texts enabling efficient searching and retrieval of the documents.

Keywords: e-content retrieval, e-learning, intelligent searching tool.

1 Introduction

With the advent of new technologies and the World Wide Web (WWW), enormous quantities of informative material are nowadays available on-line. Computers are increasingly changing from computing systems to portals, which permit to access big volumes of information. In parallel, due to the interest in reducing costs of education and stimulating people to never stop learning, e-learning applications have recently developed in educational, industrial and research institutions. However, e-learning platforms present some difficulties concerning the instructor and the student interaction. Effective mining and retrieval of the e-content is the major bottleneck of e-learning applications. The lack of metadata and classification of the e-content, force the development of powerful search engines.

The basic approaches concerning information retrieval and data mining in textual documents collections are: (1) searching using keywords or key documents, (2) exploration of the collection referring to "ome" organization or categorization of the documents, and (3) filtering of interesting documents from the incoming document stream. Keyword search systems can be automated rather easily whereas the organization of document collections has traditionally been carried out by hand. In a manual organization carried out for example, in libraries, classification schemes are defined and each document is positioned into one or several classes by a librarian.

Similarly, in the current hypertext systems the links between related documents are most often added by hand. One of the traditional methods of searching for texts that match a query is to index all the words (terms) that have appeared in the document collection. The query itself, typically a list of appropriate keywords, is compared with the term list of each document to find documents that match the query. In the existed applications, the educational content is multilingual and heterogeneous. Therefore, simple keyword queries are not capable for efficient mining of the available information. In order to bypass the aforementioned bottlenecks, we propose the use of Artificial Neural Networks (ANN). Due to their wide range of applications, ANNs have been an active research for the past decades [2]. A large variety of learning algorithms have been evolved and being employed in ANNs. A further categorization divides the network architectures into three distinct categories: feedforward, feed-backward and competitive [2]. The self-organizing maps or Kohonen's feature maps are feedforward, competitive ANN that employ a layer of input neurons and a single computational layer [7]. The neurons on the computational layer are fully connected to the input layer and are arranged on an N-dimensional lattice. In this paper, we shall use the SOM algorithm to cluster contextually similar documents into classes.

The ability of the SOM algorithm to produce spatially organized representations of the input space can be utilized in document organization, where organization refers to the representation and storage of the available data. An architecture based on the SOM algorithm that is capable of clustering documents according to their semantic similarities is the so-called WEBSOM architecture [4,5,6,7]. The WEBSOM consists of two distinct layers where the SOM algorithm is applied. The first layer is used to cluster the words found in the available training documents into semantically related collections. The second layer, which is activated after the completion of the first layer, clusters the available documents into classes that high probability contains relevant documents with respect to their semantic content. Due to that, the WEBSOM architecture regarded as a prominent candidate for document organization and retrieval. The structure of the paper is as follows. In the first section is presented the basic structure of the SOM architecture. The system architecture is presented in the second section. The software component constructed to illustrate the applicability of the proposed architecture is shown in the third section. Finally the performance of the training algorithm is illustrated in the final section.

2 Description of the SOM Algorithm

The basic Self-Organizing Map (SOM) can be visualized as a sheet-like-neural-network array (figure), the cells (or nodes) of which becomes specifically tuned to various input signal patterns or classes of patterns or classes of patterns in an orderly fashion. The learning process is competitive and supervised, meaning no teacher is needed to define the correct output (or actually the cell into which the input is mapped) for an input. In the basic version, only one map node (winner) at a time is activated corresponding to each input. The locations of the responses in the array tend to become ordered in the learning process as if some meaningful nonlinear coordinate system for the different input features were being created over the network.

Due to its competitive nature, the SOM algorithm identifies the best-matching, winning reference vector $w(k)$ (or winner for short), to a specific feature vector x_j with respect to a certain distance metric. The index s of the winning reference vector is given by:

$$s = \arg \min_{l=1}^L \| x_j - w_l(k) \|, \tag{1}$$

where $\| \cdot \|$ denotes the Euclidean distance.

A neighborhood updating, especially in the early iterations, is performed in order to achieve a global ordering of the input space onto the lattice, which is crucial for the good resolution of the map [4]. The term basic SOM will henceforth denote the on-line algorithm proposed by Kohonen [4] without any modification of speed-up techniques.

2.1 Salton’s Vector Space Model

The vector space model [8] has been widely used in the traditional IR field. Most search engines also use similarity measures based on this model to rank web documents. The model creates a space in which both documents and queries are represented by vectors. For a fixed collection of documents a dimensional vector is generated for each document and each query from sets of terms with associated weights, where m is the number of unique terms in the document collection. Then a vector similarity function, such as the inner product, can be used to compute the similarity between a document and a query. (Fig. 1)

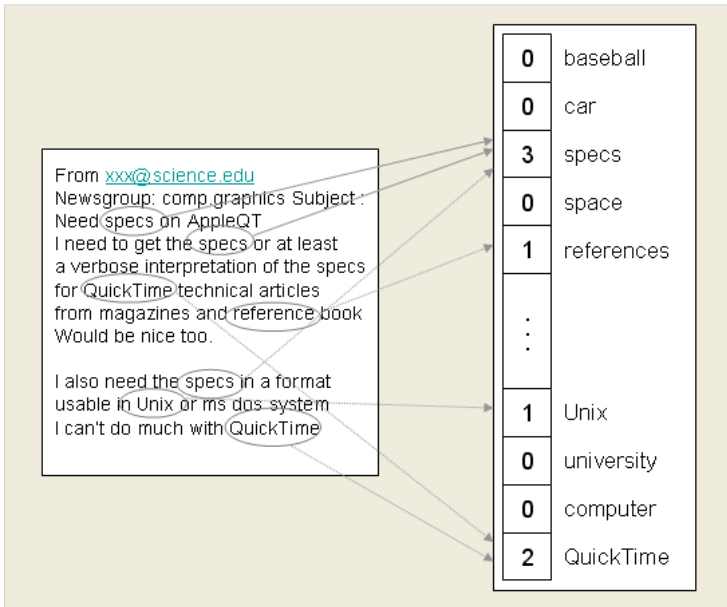


Fig. 1. The VSM. The text on the left is represented as the vector on the right. Each line of the vector is a different word of the text. Each line’s record is the frequency of the word in the text.

The main problem of the vector space model is the large vocabulary in any sizable collection of free text documents, which results in a vast dimensionality of the document vectors.

In the following sections are presented, some methods for reducing the dimensionality. These are applicable to all cases where the documents are encoded using the vector space model, i.e. as the document-by-word matrix.

2.2 Latent Semantic Indexing

Latent Semantic Indexing (LSI) is one alternative to the original vector space model. LSI tries to make account to the co-occurrence of terms in documents when encoding the documents. One way of interpreting the LSI is that it represents the j th document by the vector

$$n'_j = \sum_k n_{jk} x'_k, \quad (2)$$

where n_{jk} denotes again the number of times the word k occurs in the j th document. The x'_k is the code that the LSI forms of the k th word by investigating the co-occurrence of the words within the documents. The term-by-document matrix, a matrix where each column is the word histogram corresponding to one document, is decomposed into a set of factors (eigenvectors) using the singular-value decomposition (SVD). The factors that have the least influence on the matrix are then discarded. The motivation behind omitting the smallest factors is that they most likely consist of noise. The vector x'_i can then be formed by using only the remaining factors, whereby the dimensionality is reduced.

2.3 Random Projection

A low-dimensionality representation for documents can be obtained by a random projection of the high-dimensional representation vector into a much lower-dimensional space [3]. The benefit compared with alternative methods such as the LSI is the extremely fast computation. The accuracy of the results is still comparable.

2.4 Word Clustering

Clustering methods can be used for reducing the number of data by grouping similar items together [3]. If similar words can be clustered together, documents can be represented as histograms of words clusters rather than of individual words. Various early approaches for categorizing words have been described [3]. In languages with rigid word order, such as English, the distribution of words in the immediate context of a word contains considerable amounts of information regarding the syntactic and semantic properties of the word. The SOM has been used to cluster words based on the distributions of words in their immediate contexts [3].

2.5 SOM Computation Complexity

The computational complexity of the SOM algorithm is only linear in the number of data samples. However, the complexity depends quadratically on the number of the map units. For document maps intended for browsing the document collection, the resolution (number of map units per number of documents) should be good, since browsing is easier if there are representations of only a few, about ten documents in a map unit on the average. Hence, for such resolution, the number of map units has to be proportional to the number of documents. For very large document collections the resulting computational complexity might become problematic.

3 SOM O.D.I.S.S.E.A.S Description

The authors determine the characteristics of the intelligent search based on the nature of the educational material (e-content) that interest the users (word docs, html pages, plain text). The system has the capability of document retrieval from databases aiming at the preparation and presentation of an e-learning course. The system is capable of retrieving certain educational texts by the users according to their “physical” questions. In the following paragraph we summarize the basic elements of the system. The descriptors are exported from the text of the multimedia material and transforming of these descriptors into compound search descriptors, in suitable vectors of characteristics. The authors concretize of the non supervised learning algorithm SOM for the successful information retrieval. The concrete methodology was peered against the method of simple equation of keywords as well as the one that makes use of the simple metric resemblance in the representation space of the texts (e.g. calculation of cosine between vectors and retrieval of those nearest in the vector that represents a question) because it provides better retrieval performance and releases the user from the need of creation of complicated educational components. The big advantage of the non supervised search models is that content managers are not obliged to create huge learning material (examples of questions with the connected answers). Taking into consideration that the user cannot as an expert in neural networks training, the software search module is supposed to supply him/her with the capability to search the database intelligently, via the combination of the automatic exported characteristics and his/her own keywords.

The proposed methodology for the creation of the intelligent search system is based on the SOM algorithm, described in the previous section. In the concrete application, the SOM maps are used fro the automatic placements of the unstructured or half structured and multidimensional data of text in such a way that similar entries in the map are represented near between them. Via a learning process, the t performance is illustrated in the following section, that final map allows the direct creation of teams of terms and teams of texts so that the distances between the different data can be directly used during the search and retrieval duration. An example of a previous successful application of the SOM networks in information retrieval is the web application WebSOM [2]. This application is based on the export of descriptors of texts from different SOMs, which replace the department of pretreatment, and representation of texts, in combination with a self-organized of the

retrieved texts. It also provides the capability of a two dimension depiction of texts relative between them, so that the user has in his disposal a visual representation of the material in relevant categories. As recent researchers have shown that the functionalism of such visualization considering the help that it provides in the final user is arguable, in the concrete work, we do not use the depiction of map. On the contrary, we provide the capability of information retrieval from the database according to the content of texts and thus present the results in form of a list in a declining line of resemblance so that we decrease the difficulties faced by users.

3.1 Content Retrieval Using the SOM Algorithm

The SOM algorithm has been used to retrieve educational material from a database. The methodology of this operation follows. We export the descriptors from the text of the educational material. Responsible for this task are the designers of the database system. An automate method for transformation of the descriptors of the material as well as the compound search descriptors, in suitable vectors of characteristics is used. To perform this automatic operation we have used the VSM algorithm described previously. The VSM provides efficiency of the search results. The next step is to learn the SOM with the vectors of text characteristics. The result is the clustering of the used terms and texts in teams of relevant content. Following this step we must search for relevant texts with the use of questions. After the creation of the vector of characteristics of the question it is supplied in the entry of trained SOM network. The result is the calculation of the nearest Euclidean distance of teams of texts towards the question based on the activated neurons of the map. This team will contain texts with terms of approximate weight and consequently will present the highest affinity of content with the question. Search of relevant texts in neighbor teams is the following step. The attribute of the self-organization allows the user to search different relevant texts found in teams of neighbor neurons of map. A question that is placed to the system, formulated as word or as a combination of words, activates the processes of retrieval of texts relative with the question. The system seeks in the map of teams of terms thus neurons that correspond in the terms or in the combination of terms that exists in the question. Those texts represent in concrete neurons of the map of the teams of texts activated by the terms, are selected and thus are presented in the user. Additionally, it is possible to present texts by a concrete method of using metric resemblance between texts.

3.2 Training

After the creation of the vector maps, we perform training. The feature vectors are presented iteratively an adequate number of times to the neural network which perform clustering in an effort to build word classes containing semantically related words. This is based on empirical and theoretical observations that semantically related words have more or less the same preceding and succeeding words. The above process yields the so-called word categories map (WCM) [1]. After the computation of the document vectors the SOM method is used to cluster them. It is expected that the constructed documents classes contained contextually similar documents.

4 Experimental Results

The performance of the SOM algorithms in the proposed case study is illustrated in this section. The performance is measured using the Mean Square Error (MSE) between the reference vectors and the document vectors assigned to each neuron in the training phase. Furthermore the recall-precision performance is measured using query documents from a test set during the recall phase is used as an indirect measure of the quality of the document organization provided by the SOM algorithm.

To measure the effectiveness of a retrieval system, two widely used ratios are employed: the precision and the recall. Precision is defined as the proportion of retrieved documents that are relevant:

$$P = \frac{r}{n_2} \quad (3)$$

Recall is the proportion of relevant documents that are retrieved:

$$R = \frac{r}{n} \quad (4)$$

As the volume of the retrieved documents increases the above ratios are expected to change. The sequence of recall-precision pairs obtained yields the so-called recall-precision curve. Each query-document in the test set produces one recall-precision curve. An average over all the curves corresponding to query documents of the same topic obtained from the test set produces the average recall-precision curve. If the recall level does not equal to one, we proceed with the second best winner neuron and repeat the same procedure and so on. The comparison of the effectiveness between the retrieved documents utilizes that above-mentioned curve.

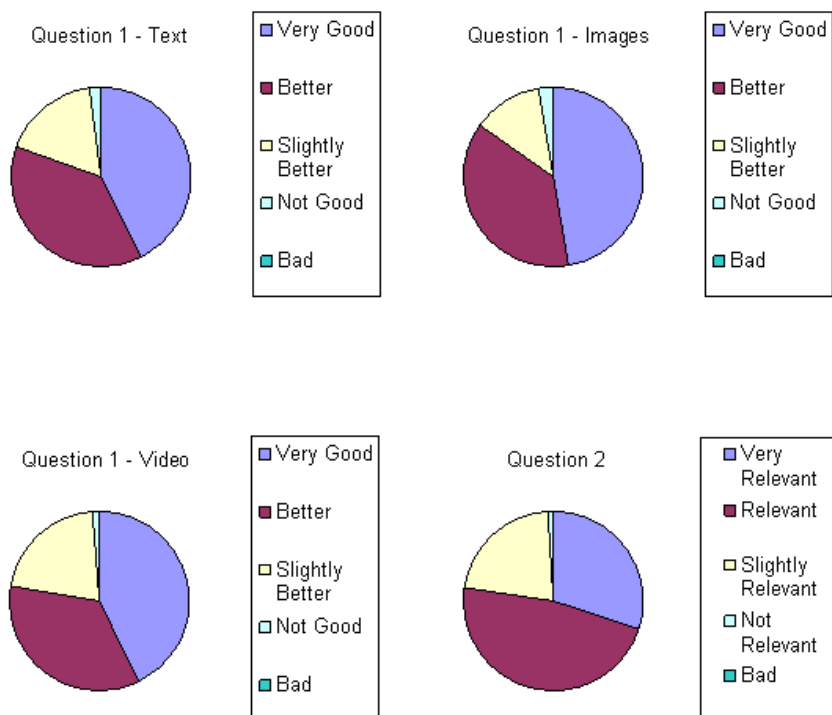
4.1 Software Component of Intelligent Search

The core of the system of the SOM algorithm and the VSM has been developed in ANSI C to ensure portability and compatibility in platforms of different type (Windows, UNIX, Linux, etc). the user interface can be designed overall round the basic and autonomous departments of the system in order to activate, through a GUI, the system operations and presentation of search results (standalone applications) via Visual Basic/C++, Delphi or Web based enabled as CGI scripting, ASP, JSP, PHP, etc. A basic issue of the final interface is the confidence estimation. The user tends to need a score of the results of its query in order to ameliorate his queering style. The following Fig. 4 presents an embodiment example of the search system into the e-learning platform O.D.I.S.S.E.A.S. [1] (Open Distance Interactive SyStem for Educational Applications) with the JSP (Java Server Pages) technology. (Table 1)

In order to prove the functionality of the proposed system, we analysed a collection of text and multimedia documents. Although the SOM algorithm can be applied only to text documents, in our collection we include multimedia documents. Multimedia

Table 1. Results obtained from the analysis of the multimedia digital library

Collection	Size (Mb)	Number of Documents	Average Document size(Mb)	Number of Lexical Form	Number of Terms
Textual	69	105	0,283	99,659	344,511
Audio (Audio Description)	90	35	2,1	35,024	140,832
Image (Image Description)	32	44	0,344	18,067	87,438
Video (Video Description)	930	16	65	3,967	23,898

**Fig. 2.** The results of the two questions

documents are analysed through the annotation that is created for each picture, video and sound. In this all the collection is available in raw text format. The test collection consists of 200 documents, 95 are multimedia and 105 are textual. The average size of the documents is 2,3 Mb and the biggest document size is 180 Mb.

4.2 Discussion

The proposed system has been evaluated by means of its usability by the e-learning users. The e-learning users are the teachers who upload new material in the database and the students who download the teaching material. In order to measure the applicability of the system, we have set two questions to the users:

Question 1: how relevant (in percentage) are the retrieved documents, in collection categories, to the query compared to ordinary search engine?

Question 2: are you satisfied with the degree of correlation of the system?

The results of the questions are depicted in Fig. 2. The only disadvantage of the proposed system is that every time a new document is uploaded in the database, the learning process must run for the new document. This process is time consuming and costly. In the future work we are planning to improve the applicability of the system concerning the automatic learning process to tune in and watch the cultural event at the same time that it's being broadcast.

5 Comparison of SOM O.D.I.S.S.E.A.S with Desktop Search Engines

SOM-O.D.I.S.S.E.A.S is focusing on the specific need of the user, which is the search of e-learning content, providing a simple and friendly user interface, in addition to the Google and Yahoo desktop search engines which provide many irrelevant add-ins and information, concerning the need of the users. In comparison to the MS search engine, SOM-O.D.I.S.S.E.A.S. provides multiple search (e.g. picture and text), in addition to the MS search that provides advanced search in categories (e.g. text or video). Also the criteria of the search are based on the "free phrase" inputted by the user and not on the name of a file or other criteria. The interface is designed following the basic rules of design, that is simple and friendly user interface which satisfies the needs of the user easily and directly, leaving out information that might distract and disorient the user from his target [14]. Thus, even a user who isn't familiar with the new technologies can simply search for e-learning content. For example, a teacher who has created a learning scenario for physics' lesson wants to search for the content he used, in order to make a scenario for the mathematics' lesson.

6 Conclusion

In this paper is presented the use of the SOM algorithm along with a training algorithm, for document retrieval. The applicability of the algorithm is illustrated in an e-learning case study. A software component has been constructed to perform intelligent search in the educational documents. The performance of the training

algorithm using the MSE measure has been presented. One of the basic issues concerning the Intelligent Systems is their ability to adjust and to be installed into any platform that requires methods of intelligent retrieval. The system was awarded after its application in various tests. Its pedagogical advantages were underlined not only by the students but also by the instructors handling and assessing the educational material. The instructors gained valuable time during their course as they could retrieve information using simple queries while students found the intelligent system necessary at their self-paced learning. Moreover, in a future expansion, the system is expected to provide reasons of the confidence estimation accompanied by the retrieved texts. That means that the user will be supplied with reasons of the certain search results and the scope of its query.

References

1. Drigas, A., Vrettaros, J., Kouremenos, S.: Open Distance Interactive System for Educational Applications O.D.I.S.S.E.A.S. In: Proceedings of the Technology and learning in Higher Education Conference, Samos (2001)
2. Honkela, J., Lagus, K., Kaski, S.: Selforganizing maps of large document collections. In: Deboeck, G., Kohonen, T. (eds.) *Visual Explorations in Finance with self-organizing Maps*, pp. 168–178. Springer, London (1998)
3. Lagus, K., Kaski, S.: Keyword selection method for characterizing text document maps. In: Proceedings of ICANN 1999, 9th international Conference on Artificial Neural Networks, vol. 1, pp. 371–376. IEE, London (1999)
4. Kohonen, T., Kaski, S., Lagus, K., Honkela, J., Paatero, V., Saarela, A.: Self organization of a massive text document collection. In: Oja, E., Kaski, S. (eds.) *Kohonen Maps*, pp. 171–182. Elsevier, Amsterdam (1999)
5. Lagus, K.: Text retrieval using self-organized document maps. Technical Report A61. Helsinki University of Technology, Laboratory of computer and information science (2000)
6. Kohonen, T.: *Self-organization and Associative Memory*, 3rd edn. Spriger, Heidelberg (1989)
7. Kohonen, T.: *Self-Organization Maps*. Springer, Heidelberg (1995)
8. Salton, A.: *Automatic Text processing*. Addition-Wesley Publishing Company, Inc., Reading (1995)
9. Kohonen, T.: Self-organization of very large documents collections. State of the art. In: Niklasson, L., Boden, M., Ziemke, T. (eds.) *Proceedings of ICANN 1998, 8th international Conference on Artificial Neural Networks*, vol. 1, pp. 65–74. IEE, London (1998)
10. Deerwester, G., Dumais, T., Furnas, W., Landauer, K., Harshman, R.: *Indexing By Latent Semantic Analysis*. Journal of the American Society of Information Science (1990)
11. Honkela, J., Kaski, S., Kohonen, T., Lagus, K.: Self-organizing maps of very large document collections: Justification for the WEBSOM method. In: Balderjahn, I., Mathar, R., Schader, M. (eds.) *Classification, Data Analysis, and Data Highways*, pp. 245–252. Springer, Berlin (1998)
12. Perfetti, R., Costantini, G.: Associative memories on BBS neural networks: a hardware-oriented learning algorithm. In: *WSEAS NNA-FSFS-EC*, pp. 109–454 (2003)
13. Triantafyllou, I., Carayannis, G.: Architectures and Techniques for Monolingual and Multilingual Information Retrieval Systems in a SOM Framework. In: *WSEAS NNA-FSFS-EC*, pp. 205–454 (2003)
14. Dix, A., Finlay, J., Abowd, G., Beale, R.: *Human Computer Interaction*, 3rd edn. Prentice Hall (2003)