# A feature mining based approach for the classification of text documents into disjoint classes

Salvador Nieto Sánchez [a], Evangelos Triantaphyllou [a,*], Donald Kraft [b]

[a] *Department of Industrial and Manufacturing Systems Engineering, 3128 CEBA Building, Louisiana State University, Baton Rouge, LA 70803, USA*
[b] *Department of Computer Science, 286 Coates Hall, Louisiana State University, Baton Rouge, LA 70803, USA*

## Abstract

11    This paper proposes a new approach for classifying text documents into two disjoint classes. The new
12  approach is based on extracting patterns, in the form of two logical expressions, which are defined on
13  various features (indexing terms) of the documents. The pattern extraction is aimed at providing de-
14  scriptions (in the form of two logical expressions) of the two classes of positive and negative examples. This
15  is achieved by means of a data mining approach, called One Clause At a Time (OCAT), which is based on
16  mathematical logic. The application of a logic-based approach to text document classification is critical
17  when one wishes to be able to justify why a particular document has been assigned to one class versus the
18  other class. This situation occurs, for instance, in declassifying documents that have been previously
19  considered important to national security and thus are currently being kept as secret. Some computational
20  experiments have investigated the effectiveness of the OCAT-based approach and compared it to the well-
21  known vector space model (VSM). These tests also have investigated finding the best indexing terms that
22  could be used in making these classification decisions. The results of these computational experiments on a
23  sample of 2897 text documents from the TIPSTER collection indicate that the first approach has many
24  advantages over the VSM approach for solving this type of text document classification problem. More-
25  over, a guided strategy for the OCAT-based approach is presented for deciding which document one needs
26  to consider next while building the training example sets. © 2002 Elsevier Science Ltd. All rights reserved.

27  *Keywords:* Document classification; Document indexing; Vector space model; Data mining; One Clause At a Time
28  (OCAT) algorithm; Machine learning

* Corresponding author. Tel.: +1-255-388-5372; fax: +1-255-388-5990.
  *E-mail addresses:* snieto@gi.com (S. Nieto Sánchez), trianta@lsu.eduhttp://www.imse.lsu.edu/vangelis/ (E. Trianta-phyllou), kraft@bit.csc.lsu.edu (D. Kraft).

## 1. Introduction and background information

This paper investigates the problem of classifying text documents into two disjoint classes. Two sample sets of training examples (text documents) are assumed to be available. An approach is developed that uses indexing terms to form logical expressions (patterns) that next are used to classify unseen text documents. This is a typical case of supervised "crisp" classification.

A typical application of this type of classification problem occurs in the declassification process of vast amounts of documents originally produced by the US Federal Government. For reasons of national security, today there are huge numbers of documents that remain classified as secret. These documents are being kept in secured places because once they were considered to be important to national security. However, high maintenance costs and new laws dictate that these documents should be re-evaluated, and the ones that are not critical any more should become public. Thus, such documents need to be (roughly speaking) classified into the two disjoint categories of "*secret*" and "*non-secret*". In this context, when a document becomes "*non-secret*" after being "*secret*", it is termed as "*declassified*". In reality, documents have been classified into multiple levels of criticality to the national security, but in this study we will consider only two classes as described above. It should also be stated here that once a document becomes public (i.e., it has been declassified), then there is no way to make it secret again (especially now with the proliferation of the Internet).

In order to highlight the complexity of this kind of classification problem, consider the evaluation of the following three illustrative sentences: "*An atomic test is to be conducted at Site X*", "*An atomic test is to be conducted at 1:00 p.m.*", and "*An atomic test is to be conducted at Site X at 1:00 p.m.*" which come from three hypothetical documents, *A*, *B*, and *C*, respectively. According to (DynMeridian, 1996) and (DOE, 1995), only document *C* is both specific and sensitive and should not be declassified and instead should continue to be kept secret. The reason for this DOE (US Department of Energy) classification rule is because document *C* includes a sentence with specific reference to the "*place and time*" of an "*atomic test*". On the other hand, documents *A* and *B* can be declassified (assuming that the rest of their contents is not critical) and become available to the general public. In this illustrative example some key text features that can be used to characterize the two classes are references to "*place*", "*time*" and "*atomic test*".

Traditionally, this declassification process is carried out by employing vast numbers of human experts. However, the sheer amount of documents under consideration can make this process extremely ineffective and inefficient. Although there are guidelines of how to declassify secret documents, directly computerizing the human effort would require developing sophisticated parsers. Such parsers would have to analyze syntactically a document and then determine which, if any, guideline is applicable. The poor quality of the documents (many of which are decades old) and the complexity of the declassification guidelines could make such an approach too risky to national security. Thus, a reasonable alternative is to seek to employ machine-learning techniques. More specifically, techniques that use "learning from examples" approaches might be appropriate. Thus, this study is centered on the following three inter-related research problems:

1. Employ data mining techniques for extracting (mining) from two sets of examples text features that could be used to correctly group documents into two disjoint classes.

2. Use such features to form logical expressions (patterns) that could explain how the training examples are grouped together, and accurately classify unseen documents.

72    3. When considering a guided learning strategy for extracting the logical expressions, determine
73      the next training document to include in the sets with the training examples so that accurate
74      logical expressions are extracted as quickly as possible.

75    Since being able to justify this kind of classification decisions is important (given the severity of
76 wrongly releasing a sensitive document to the public), methods that do not clearly allow for an
77 explanation of the decision making process are not appealing here. Therefore, an impetus for this
78 research is to seek to develop an approach that is based on mathematical logic, versus approaches
79 that do not provide satisfactory explanation capabilities.

80    Traditional text classification and information retrieval (IR) techniques may have some limi-
81 tations in solving this problem because they group documents that share a similar content. The
82 prime example of such techniques is the vector space model (VSM) (Salton, 1989), which according
83 to the literature (Buckley & Salton, 1995; Shaw, 1995) is the most effective methodology for this
84 type of classification. The limitation of this technique is that it is based on similarity measures and
85 thus it may not be able to distinguish between the previous three illustrative sentences in terms of
86 the critical classification issues despite their all having similar contents. Other techniques, such as
87 fuzzy set approaches (FSAs), neural networks (NNs), nearest neighbor, and computational se-
88 mantic analysis (SA), have limitations in addressing these research problems, either because of
89 their time complexity or because the resulting sizes of their outputs are still unacceptable and do
90 not possess satisfactory explanatory capabilities (Chen, 1996; Scholtes, 1993).

91    An alternative approach to address the present research problems is the *One Clause At a Time*
92 (*or OCAT*) algorithm (Triantaphyllou, 1994; Triantaphyllou, Soyster, & Kumara, 1994). This is a
93 new data mining approach based on mathematical logic. This approach extracts (mines) key
94 features from the training examples and next uses them to construct logical expressions (patterns)
95 that can be used in classifying the training examples into the two original disjoint classes. These
96 logical expressions can also easily be transformed into the IF-THEN type of decision rules. The
97 OCAT approach applies to examples that can be represented by binary vectors (although attri-
98 butes with continuous values can be transformed into ones with binary values (Triantaphyllou,
99 2001). However, this is not a limitation because it is the mere presence or absence of certain key
100 words that can cause a document to be grouped in one class or another. On the other hand, the
101 typical document classification done by traditional IR systems uses term frequencies (which are
102 continuous numbers usually normalized in the interval $[0, 1]$) of keywords to group together
103 documents of seemingly similar context.

104    This paper is organized as follows. Section 2 presents an overview of the document indexing
105 process. Section 3 introduces the OCAT algorithm. Section 4 presents an overview of the VSM
106 algorithm. Section 5 presents an overview of the guided learning approach (GLA). Section 6
107 describes the methodology of this investigation. Section 7 presents and summarizes the results.
108 Finally the paper ends with some concluding remarks.

109 **2. A brief overview of the document clustering process**

110    The traditional process for automatic clustering of text documents results in a grouping of
111 documents with similar content into meaningful groups in order to facilitate their storage and
112 retrieval (see, for example, Salton, 1989). This is a four-step process as follows. In the first step a

4              *S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*

113  computerized system compiles a list of the unique words that co-occur in a sample of the doc-
114  uments from various classes (see, for example, Cleveland & Cleveland, 1983; Salton, 1989). In the
115  second step, the co-occurring frequency of these words is analyzed and the best set of indexing
116  terms is extracted. Usually, indexing terms (also known as *keywords* or *content descriptors*) are
117  selected among the words with moderate frequency. The most *common* and the most *rare* words
118  (i.e., the most frequent and infrequent words, respectively) are discarded as keywords because
119  they convey little lexical meaning (see, for example, Cleveland & Cleveland, 1983; Fox, 1990;
120  Luhn, 1957, 1958; Meadow, 1992; Salton, 1968; Zipff, 1949). Some examples of common words
121  are: "a", "an", "and", and "the" (Fox, 1990); rare words depend on a document's subject domain
122  (Meadow, 1992).
123      In the third step, a document is indexed by affixing it with the set of keywords that only occur in
124  its text. According to Cleveland & Cleveland (1983), "this assignment is correct because authors
125  usually repeat words that conform with the document's subject". A list (vector) of keywords
126  represents the content of a document and usually it is referred to as a *document representation* or
127  *surrogate*. An example of such a surrogate is the list of the seven words or phrases: {"*Document
128  classification*", "*document indexing*", "*vector space model*", "*data mining*", "*OCAT algorithm*",
129  and "*machine learning*"}. This surrogate could be used to represent the content of this paper,
130  which is composed of thousands of words, symbols, and numbers. Hence, the goal of the third
131  step is to construct a surrogate for representing the content of each document.
132      An advantage of using such surrogates is that they can be further simplified by expressing them
133  as numerical vectors (Salton, 1989). One way to construct such vectors is by expressing their
134  elements as binary values to indicate the presence (denoted by 1) or absence (denoted by 0) of
135  certain keywords in a document. For instance, the vector's element $w_{ij} = 1$ (or 0) expresses the
136  presence (or absence) of keyword $T_i$ ($i = 1, 2, 3, \ldots, t$) in document $D_j$ ($j = 1, 2, 3, \ldots, N$). Thus,
137  the surrogate $D_j = [0\,1\,1\,1\,1\,0]$ of six binary elements indicates the presence of keywords (terms) $T_2$,
138  $T_3$, $T_4$, and $T_5$ and the absence of keywords (terms) $T_1$ and $T_6$ in $D_j$.
139      Another way to construct these numerical vectors is by expressing (i.e., weighting) their ele-
140  ments using real values from the range [0, 1]. In this case, the value of an element $w_{ij}$ indicates the
141  relative occurrence frequency of a keyword within a document. For instance, a hypothetical
142  surrogate such as $D_j = [0.00 \quad 1.00 \quad 0.10 \quad 0.75 \quad 0.90 \quad 1.00]$ may indicate that term $T_3$ occurs
143  little, terms $T_4$ and $T_5$ occur moderately, and terms $T_2$ and $T_6$ occur with high frequency in $D_j$. In
144  the remainder of this paper, however, only binary surrogates will be considered. As stated above
145  the reason for considering binary vectors as surrogates is because the mere presence or absence of
146  a keyword (or some pattern of keywords) may be detrimental in assigning a document to one of
147  the two disjoint classes considered in this paper.
148      In the last step of the (traditional) classification process, documents sharing similar keywords
149  (i.e., content) are grouped together. This classification follows from the pairwise comparison of all
150  the surrogates (Salton, 1989).

## 3. An overview of the OCAT algorithm

152      The OCAT algorithm (Triantaphyllou, 1994; Triantaphyllou et al., 1994) is an inductive
153  learning (data mining) approach for the classification of examples (documents in this study) into

5

154 two disjoint classes. Each example is represented by a binary vector, although it can be generalized
155 into vectors defined on continuous variables (Triantaphyllou, 2001). The $i$th element of such a
156 vector represents the presence (1) or absence (0) of a key characteristic (also called a feature,
157 variable, atom, or attribute) pertinent to the phenomenon under study. For this study, such bi-
158 nary features represent the presence or absence of the index terms (keywords) in the document
159 surrogates. It is a good idea for the analyst to first try to define these examples by using as many
160 features as possible.

161 The OCAT algorithm systematically identifies ("mines") a small set of features while simul-
162 taneously constructing a logical expression of small size defined on these features that can be used
163 to group together the training examples into the two original disjoint classes. This logical ex-
164 pression is a Boolean function that evaluates to true when it is given training examples from one
165 of the two disjoint classes (the "positive" class) and false when it is given training examples from
166 the other class (the "negative" class). The assignment of the terms "positive" and "negative" is
167 arbitrary. Hopefully, if the training examples are representative enough, then this logical ex-
168 pression can be used to accurately classify unseen examples. Moreover, this logical expression can
169 be used to extract new knowledge pertinent to the system under study. The extracted Boolean
170 function is in conjunctive normal form (CNF) or in disjunctive normal form (DNF). Fig. 1 il-
171 lustrates this algorithm for the CNF case.

172 From Fig. 1 it becomes evident that the OCAT algorithm is greedy in nature. This allows it to
173 return logical expressions of small size. That is, it attempts to return a logical expression that is
174 comprised of a small number of CNF (or DNF) clauses (Triantaphyllou & Soyster, 1996a). This
175 algorithm is greedy in the sense that in the first iteration, it forms a clause that for the CNF case
176 accepts (i.e., it evaluates to 1) all the examples in one of the classes ($E^+$) and rejects (i.e., it
177 evaluates to 0) as many examples in the other class ($E^-$) as possible. Similarly, in the second it-
178 eration it forms another clause that again accepts all the examples in $E^+$ and rejects as many

**Input**:       $E^+$ and $E^-$ are the two disjoint sets with the "positive" and "negative" training
               examples, respectively.
**Output**:      A Boolean function in CNF (for this implementation) form.

        **Begin**
          $i = 0$; $C = \varnothing$;   /* initializations */
          **do while** ($E^- \neq \varnothing$ )
*Step 1*:     Set $i \leftarrow i + 1$; /* where $i$ indicates the $i$-th iteration */
*Step 2*:     Find a clause $C_i$ which accepts all members of $E^+$ while it rejects as many
               members of $E^-$ as possible;
*Step 3*:     Let $E^-(C_i)$ be the set of the members of $E^-$ which are rejected
               by the CNF clause $C_i$;
*Step 4*:     Let $C \leftarrow C \chi C_i$;
*Step 5*:     Let $E^- \leftarrow E^- - E^-(C_i)$;
          **repeat**;
        **End;**

Fig. 1. The OCAT algorithm for CNF expressions.

6                 *S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*

179 examples as possible in $E^-$ that were accepted by the previous clause(s). The algorithm repeats this
180 process until all the examples in $E^-$ are rejected by the generated sequence of clauses.
181     Triantaphyllou et al. (1994) and Triantaphyllou (1994) have implemented a branch-and-bound
182 (B&B) approach to construct a set of clauses of minimal cardinality (i.e., a logical expression
183 comprised by a minimum number of clauses) to solve the problem in *Step 2* of the OCAT al-
184 gorithm in Fig. 1. This approach, however, is limited to problems of small to medium size because
185 of the extensive CPU times it takes to find an optimal or next to optimal set of clauses (i.e., a set of
186 clauses of minimal cardinality). In Triantaphyllou & Soyster (1995) the authors present a simple
187 transformation approach that can be used to extract DNF functions from algorithms that infer
188 CNF functions and vice versa.
189     More recently, Deshpande & Triantaphyllou (1998) have implemented a faster heuristic of
190 quadratic time complexity as it is highlighted in Fig. 2. The version of the OCAT approach de-
191 picted in Fig. 2 returns a CNF expression (Boolean function) which accepts all of the positive
192 examples while rejecting all of the negative examples. This algorithm can also be enhanced with a
193 randomized approach which essentially solves this problem many times (in a randomized fashion)
194 and at the end it selects the smallest (in terms of the number of clauses in the Boolean expression)
195 of the extracted functions. This solution allows for larger classification problems and delivers a set
196 of (CNF or DNF) clauses of small size (as opposed to the minimal or near to minimal size sets
197 delivered by the B&B approach).
198     The notation $POS(A_j)$ and $NEG(A_j)$ in Fig. 2 denotes the number of positive and negative
199 examples that will be accepted, respectively, if feature (atom) $A_j$ is introduced in the current CNF
200 clause. That is, $POS(A_j) = |E^+(A_j)|$ and $NEG(A_j) = |E^-(A_j)|$. Please note that Steps 1 and 2 in
201 Fig. 2 also consider the negations (denoted as $\bar{A}_j$) of the features $A_j$. The ratio $POS(A_j)/NEG(A_j)$
202 in Steps 1 and 2 is used to quickly identify features that can form a clause that would accept all the

**Input**:      $E^+$ and $E^-$ are the two training sets as before defined on the binary atoms $A_j$
            (for $j = 1, 2, 3, ..., t$).
**Output**:   A Boolean function in CNF which accepts all positive and rejects all negative training
            examples.

Set $k = 1$ ;  Set $C = \varnothing$ ; /* initializations */
**do while** ($E^-$  $\varnothing$)
        Let $E^+$ be the original set of positive training examples ;
        Set $C_k = \varnothing$;  /* initialize the current clause */
        **do while** ($E^+$  $\varnothing$)
         *Step 1*: Calculate the $POS(A_j) / NEG(A_j)$ ratio for all features  $A_j$ (and their negations);
         *Step 2*: Choose a new feature $A_j$ according to $max[POS(A_j) / NEG(A_j)]$ value.
         *Step 3*: Let $C_k \leftarrow C_k \vee A_j$ ;
         *Step 4*: Let  $E^+(A_j)$ be the set of members of $E^+$ which are accepted when $A_j$ is
                included in the current clause $C_k$ ;
         *Step 5*: Let $E^+ \leftarrow E^+ - E^+(A_j)$ ;
        **repeat**;
        *Step 6*.   Let $E^-(C_k)$ be the set of members of $E^-$ which are rejected by $C_k$ ;
        *Step 7*:   Let $E^- \leftarrow E^- - E^-(C_k)$ ;  Let $k \leftarrow k + 1$ ;  Let $C \leftarrow C \wedge C_k$ ;
**repeat**;

Fig. 2. A fast heuristic for forming CNF clauses for the OCAT approach.

203 positive examples while rejecting many of the remaining negative ones (for the CNF case). By
204 selecting a feature that maximizes this ratio, it is likely to have a feature that has a large $\text{POS}(A_j)$
205 value and a low $\text{NEG}(A_j)$ value. If the value of $\text{NEG}(A_j)$ is equal to zero, then a large value is
206 assigned to the ratio in Step 1 (in Fig. 2) in order to make the choice of including the feature
207 (atom) $A_j$ very appealing.

208    In order to help illustrate the previous issues consider the two sets of training examples depicted
209 in Fig. 3. The set with the positive examples is comprised of four examples, while the set of the
210 negative examples is comprised of six examples. These examples (document surrogates in our
211 context) are defined on four binary features (i.e., index terms or atoms) or their negations. A value
212 of 1 indicates the presence of the corresponding index term, while a value of 0 indicates the ab-
213 sence of the index term.

214    In the experiments described in this paper, the OCAT approach (as depicted in Fig. 1) employs
215 the fast heuristic shown in Fig. 2. However, in other implementations a B&B approach (Tri-
216 antaphyllou, 1994) may be used in Step 2. When the fast heuristic in Fig. 2 is employed, the first
217 CNF clause that is derived is $(A_2 \vee A_4)$. To follow this, observe that the ratio with the maximum
218 value is given by $\text{POS}(A_2)/\text{NEG}(A_2) = 2/2 = 1.00$ (this is a tie with the ratio that corresponds to
219 $A_4$). When feature $A_2$ (which is arbitrarily selected over $A_4$) is included in the first clause (which
220 initially is nil), then the first and the second positive examples will be accepted by that clause. The
221 next iteration of the fast heuristic will consider the same negative examples, but the updated set of
222 positive examples consists of the remaining positive examples (i.e., the third and fourth). This will
223 result in the next feature to be selected being $A_4$. Now, the CNF clause that is comprised of these
224 two features (i.e., $A_2$ and $A_4$) will accept all the positive examples in $E^+$ while rejecting the first,
225 fourth and fifth examples from the set of the negative examples in $E^-$.

226    The next iteration of the OCAT approach (in Fig. 1) will consider the original four positive
227 examples, but now the set of the negative examples consist of the ones not rejected so far. That is,
228 the second, third, and sixth negative examples, in terms of the original set $E^-$. Working as above,
229 it can be seen that the next application of the fast heuristic will return the clause : $(\bar{A}_2 \vee \bar{A}_3)$. The
230 loop in Fig. 1 needs to be repeated once more, and a third (and final) clause is derived. That clause
231 is $(A_1 \vee A_3 \vee \bar{A}_4)$. In other words, the logical expression (in CNF) which is derived from the
232 training examples depicted in Fig. 3 is as follows:

$$(A_2 \vee A_4) \wedge (\bar{A}_2 \vee \bar{A}_3) \wedge (A_1 \vee A_3 \vee \bar{A}_4). \tag{1}$$

234    A fundamental property of expression (1) is that it accepts (i.e., evaluates to 1) all the examples
235 in $E^+$, while it rejects (i.e., evaluates to 0) all the examples in $E^-$. However, since such an ex-

$$E^+ = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \qquad E^- = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

Fig. 3. A set of four positive and a set of six negative examples.

*S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*

236 pression is usually constructed from a relatively small collection of training examples, it is possible
237 that the expression to be inaccurate when it classifies unseen examples. The error may occur if
238 either the unseen example is positive, and the expression rejects it, or if the example is negative and
239 the expression accepts it. More details on the OCAT approach can be found at: http://www.im-
240 se.lsu.edu/vangelis.

## 4. Improving the classification accuracy of the OCAT algorithm by using two logical expressions

242 Let us consider generating a second logical expression for classifying unseen examples. Please
243 recall that the first expression is derived by treating the first set of training examples as positive
244 and the second as the negative examples. However, the second expression is derived by treating
245 the second set of training examples as positive and the first as the negative training examples. For
246 instance, Fig. 4 depicts the same training examples as the ones in Fig. 3, but now they have reverse
247 roles.

248 When the OCAT approach is applied on the new inference problem, the following expression
249 (2) is derived:

$$(A_3 \vee \bar{A}_2) \wedge (\bar{A}_4 \vee A_2 \vee \bar{A}_1) \wedge (A_1 \vee \bar{A}_3). \tag{2}$$

251 As with expression (1), a property of the corresponding Boolean function $f(x) = (A_3 \vee \bar{A}_2) \wedge$
252 $(\bar{A}_4 \vee A_2 \vee \bar{A}_1) \wedge (A_1 \vee \bar{A}_3)$ is to accept (i.e., to evaluate to 1) the former negative examples and to
253 reject (i.e., to evaluate to 0) the former positive examples. For convenience, following the setting
254 of the examples in Figs. 3 and 4, expression (1) will be called the *positive rule* (denoted as $R^+$) while
255 expression (2) the *negative rule* (denoted as $R^-$).

256 The disadvantage of using only one rule (logical expression) can be overcome by considering
257 the combined decisions of $R^+$ and $R^-$ when classifying an unseen example $e$. If $e$ is a positive
258 example it will be denoted as $e^+$, while if it is a negative example it will be denoted as $e^-$. Under
259 this setting, the classification of $e$ can only be:

260 1. *Correct* if and only if:
261     (a) $R^+(e^+) = 1$ and $R^-(e^+) = 0$;
262     (b) $R^+(e^-) = 0$ and $R^-(e^-) = 1$.

263 2. *Incorrect* if and only if:
264     (c) $R^+(e^+) = 0$ and $R^-(e^+) = 1$;
265     (d) $R^+(e^-) = 1$ and $R^-(e^-) = 0$.

$$E^+ = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 \end{bmatrix} \qquad E^- = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

Fig. 4. The training example sets in reverse roles.

266  3. *Undecided* if and only if:
267       (e) $R^+(e^+) = 1$ and $R^-(e^+) = 1$;
268       (f) $R^+(e^-) = 1$ and $R^-(e^-) = 1$;
269       (g) $R^+(e^+) = 0$ and $R^-(e^+) = 0$;
270       (h) $R^+(e^-) = 0$ and $R^-(e^-) = 0$.
271  Cases (a) and (b) are called ''correct'' classifications because both rules perform according to
272  the desired properties described above. However, as indicated above it is possible that the rules
273  could incorrectly classify an example (cases (c) and (d)). Or the rules could simultaneously accept
274  (cases (e) and (f)) or reject (cases (g) and (h)) the example. Cases (e)–(h) are called ''undecided''
275  because one of the rules does not possess enough classification knowledge, and thus such a rule
276  must be reconstructed. Therefore, ''undecided'' situations open the path to improve the accuracy
277  of a classification system. This paper exploits the presence of ''undecided'' situations in order to
278  guide the reconstruction of the rule that triggered an erroneous classification decision.

279  **5. An overview of the vector space model**

280  The VSM is a mathematical model of an IR system that can also be used for the classification of
281  text documents (Salton, 1989; Salton & Wong, 1975). It is often used as a benchmarking method
282  when dealing with document retrieval and classification related problems. Fig. 5 illustrates a
283  typical three-step strategy of the VSM approach to clustering.
284  To address *Step 1* Salton (1989) indicates that a suitable measure for pairwise comparing any
285  two surrogates $X$ and $Y$ is the cosine coefficient (CC) as defined in Eq. (3) (other similarity
286  measures are listed in Salton (1989, Chapter 10):

$$CC = \frac{|X \cap Y|}{|X|^{1/2} \cdot |Y|^{1/2}}. \tag{3}$$

288  In this formula, $X = (x_1, x_2, x_3, \ldots, x_t)$ and $Y = (y_1, y_2, y_3, \ldots, y_t)$, where $x_i$ indicates the presence
289  (1) or absence (0) of the $i$th indexing term in $X$ and similarly for $y_i$ in respect to $Y$. Moreover,
290  $|X| = |Y| = t$ is the number of indexing terms, and $|X \cap Y|$ is the number of indexing terms ap-
291  pearing simultaneously in $X$ and $Y$. To be consistent with the utilization of binary surrogates,
292  formula (3) provides the CC expression for the case of Boolean vectors. This coefficient measures
293  the angle between two surrogates (Boolean vectors) in a Cartesian plane. Salton (1989) indicates
294  that ''the magnitude of this angle can be used to measure the *similarity* between any two docu-

**Input**:          A sample of surrogates
**Output**:          Clusters of documents and clusters' centroids

*Step 1*:  Compute the pairwise similarity coefficients among all surrogates in the sample
*Step 2*  Cluster documents with sufficiently large pairwise similarities
*Step 3*  Compute the centroids of the clusters

Fig. 5. The VSM approach.

295 ments''. This statement is based on the observation that two surrogates are identical, if and only if
296 the angle between them is equal to 0°.

297    In *Step 2* the VSM clusters together documents that share a similar content based on their
298 surrogates. According to Salton (1989), any clustering technique can be used to group documents
299 with similar surrogates. A collection of clustering techniques is given in Anderberg (1973), Ald-
300 enderfer (1984), Späth (1985), and Van Rijsbergen (1979). However, it is important to mention
301 here that with any of these techniques, the number of generated classes is always a function of
302 some predefined parameters. This is in contrast with the requirements of our problem here in
303 which the number of classes is exactly equal to two. When the VSM works under a predefined
304 number of classes, it is said to perform a pseudo-classification. In this study the training examples
305 are already grouped into two (disjoint) classes. Thus, the VSM is applied on the examples
306 (documents) in each class and the corresponding centroids are derived. Hence, we will continue to
307 use this kind of pseudo-classification.

308    To address *Step 3* Salton (1989), Salton & Wong (1975), and Van Rijsbergen (1979) suggest the
309 computation of a class centroid to be done as follows. Let $w_{rj}$ $(j = 1, 2, 3, \ldots, t)$ be the $j$th element
310 of the centroid for class $C_r$ which contains $q$ documents. Also, the surrogate for document $D_i$ is
311 defined as $\{D_{ij}\}$. Then, $w_{rj}$ is computed as follows:

$$w_{rj} = (1/q) \sum_{i=1}^{q} D_{ij} \quad \text{for } j = 1, 2, 3, \ldots, t. \tag{4}$$

313 That is, the centroid for class $C_r$ is also a surrogate (also known as the "average" document)
314 defined on $t$ keywords.

315    Finally, the VSM classifies a new document by comparing (i.e., computing the CC) its surrogate
316 against the centroids that were created in *Step 3*. A new document will be placed in the class for
317 which the CC value is maximum.

318    In the tests to be described later in this paper, the VSM is applied on the documents (training
319 examples) available for each class. In this way, the centroid of each one of the two classes is
320 derived. For instance, consider the training examples depicted in Figs. 3 and 4. The VSM is now
321 applied on these data. The centroids in expression (5) have been constructed from the data in Fig.
322 3 and the centroids in expression (6) from the data in Fig. 4. Obviously, the centroids for the
323 second set are in reverse order of those for the first set of data.

$$\begin{aligned} C_+ &= [1/2, 1/2, 1/4, 1/2], \\ C_- &= [2/3, 1/3, 1/2, 1/3], \end{aligned} \tag{5}$$

$$\begin{aligned} \mathbb{C}_+ &= [2/3, 1/3, 1/2, 1/3], \\ \mathbb{C}_- &= [1/2, 1/2, 1/4, 1/2]. \end{aligned} \tag{6}$$

326    The notation $C_+$ and $C_-$ stands for the centroids for the data in Fig. 3, and $\mathbb{C}_+$ and $\mathbb{C}_-$ stand for
327 the centroids for the data in Fig. 4, respectively. In order to match the names of the positive and
328 negative rules described for the OCAT algorithm, the two centroids for the data in Fig. 3 will be
329 called the *positive centroids* while the centroids for the data in Fig. 4 will be called the *negative*
330 *centroids*. As with the OCAT algorithm, the utilization of two sets of centroids has been inves-

331 tigated in order to tackle the new classification problem by using the VSM as new examples
332 become available.

## 6. Guided learning for the classification of text documents

334    The central idea of the GLA can be illustrated as follows. Suppose that the collection to be
335 classified contains millions of documents. Also, suppose that an oracle (i.e., an expert classifier) is
336 queried in order to classify a small sample of examples (documents) into classes $E^+$ and $E^-$. Next,
337 suppose that the OCAT algorithm is used to construct the positive and negative rules, such as was
338 the case with expressions (1) and (2). As indicated earlier, these rules may be inaccurate when
339 classifying examples not included in the training set, and therefore they will result in one of the
340 classification outputs provided in cases (a)–(g), as described earlier. One way to improve the
341 classification accuracy of these rules is to add one more document to the training set (either in $E^+$
342 or $E^-$) and have them reconstructed. Therefore, the question GLA attempts to answer is: *What is*
343 *the next document to be inspected by the expert classifier so that the classification performance of the*
344 *derived rules can be improved as fast as possible*?
345    One way to provide the expert with this document is to randomly select one from the remaining
346 unclassified documents. We call this the RANDOM input learning strategy. A drawback of this
347 strategy may occur if the oracle and *incumbent* rules frequently classify a document in the same
348 class. If this occurs frequently, then the utilization of the oracle and the addition of the example to
349 the training set is of no benefit. An alternative and more efficient way to provide the expert with a
350 document is to select one in an ''undecided'' situation. This strategy (in a general form) was first
351 introduced in Triantaphyllou & Soyster (1996b). This approach appears to be a more efficient way
352 of selecting the document because an ''undecided'' situation implies that one of the rules mis-
353 classified the document. Therefore, the expert's verdict will not only guide the reconstruction of
354 the rule that triggered the misclassification, but it may also improve the learning rate of the two
355 rules. We call this the GUIDED input learning strategy. An incremental learning version of the
356 OCAT approach is described in Nieto Sanchez, Triantaphyllou, Chen, & Liao (2001).

## 7. Experimental data

358    In order to determine the classification performance of the OCAT approach in addressing this
359 new problem, the OCAT's classification accuracy was compared with that of the VSM. Both
360 approaches were tested under three experimental settings:
361 1. a Leave-One-Out Cross-Validation (or CV) (also known as the Round-Robin test);
362 2. a 30/30 Cross-Validation (or 30CV), where 30 stands for the number of training documents in
363    each class; and
364 3. in an experimental setting in which the OCAT algorithm was studied under a random and a
365    guided learning strategy.
366 These will be defined below. This multiple testing strategy was selected in order to gain a more
367 comprehensive assessment of the effectiveness of the various methods.

368  For these tests, a sample of 2897 documents was randomly selected from four document classes
369  of the TIPSTER collection (Harman, 1995; Voorhees, 1998). The previous numbers of documents
370  in each class were determined based on memory limitations on the computing platform used (an
371  IBM Pentium II PC running Windows 95). The TIPSTER collection is a standard data set for
372  experimentation with IR systems. The four document classes were as follows:
373  1. Department of Energy (DOE) documents,
374  2. Wall Street Journal (WSJ) documents,
375  3. Associated Press (AP) documents, and
376  4. ZIPFF class documents.
377  We chose documents from this collection because for security reasons we did not have access to
378  actual secret DOE documents.
379  Table 1 shows the number of documents that were used in the experimentation. These docu-
380  ments were randomly extracted from the four classes of the TIPSTER collection.
381  We simulated two mutually exclusive classes by forming the following five class-pairs: (DOE vs.
382  AP), (DOE vs. WSJ), (DOE vs. ZIPFF), (AP vs. WSJ), and (WSJ vs. ZIPFF). These five class-
383  pairs were randomly selected from all possible class-pairs combinations. To comply with the
384  notation presented in the previous sections, the first class of each class-pair was denoted as $E^+$,
385  while the second class was denoted as $E^-$. Thus, we try to find a Boolean expression to classify a
386  document surrogate into the proper TIPSTER class.
387  Table 2 shows the average number of keywords that were extracted from the five class-pairs
388  mentioned above. The data in this table can be interpreted as follows. For the class-pair (DOE vs.
389  AP), the average number of keywords used in all the experiments was 511 under the CV validation
390  and 803 under the 30CV validation. A similar interpretation applies to the data in the other
391  columns.
392  It should be stated at this point that we used a number of alternative indexing terms. Besides
393  single words, we also used sequences of two words at a time, sequences of three words at a time,
394  and sequences of four words at a time. However, some pilot studies indicated that the best results

Table 1
Number of documents randomly extracted from each class

| Class: | DOE | AP | WSJ | ZIPFF | Total |
|---|---|---|---|---|---|
| Number of documents: | 1407 | 336 | 624 | 530 | 2897 |

DOE, AP, WSJ, and ZIPFF stand for Department of Energy, Associated Press, and the Wall Street Journal, respectively; ZIPFF is a collection of technical documents of various topics.

Table 2
Average number of indexing words used in each experiment

| Type of experiment | DOE vs. AP | DOE vs. WSJ | DOE vs. ZIPFF | AP vs. WSJ | WSJ vs. ZIPFF |
|---|---|---|---|---|---|
| CV | 511 | 605 | 479 | 448 | 501 |
| 30CV | 803 | 911 | 890 | 814 | 811 |

In order to keep the size reasonable for our computing environment, only the first hundred words from each document were considered. Stop words were always removed.

395  would be derivable by using as indexing terms single words only. Thus, the atoms (binary vari-
396  ables) in the derived logical expressions are single keywords and not sequences of them.

397  **8. Testing methodology**

398       This section first summarizes the methodology for the Leave-One-Out Cross-Validation and
399  the 30/30 Cross-Validation. These two alternative testing methods have been employed in order to
400  gain a better understanding of the various procedures used to classify text documents. The same
401  section also presents the statistical tests employed to determine the relative performance of the
402  VSM and the OCAT algorithm. This section ends with the methodology for the GLA.

403  *8.1. The Leave-One-Out Cross-Validation*

404       The cross-validation (CV) testing was implemented on samples of 60 documents as follows.
405  First, 30 documents from each class were randomly selected. Please note that the size 60 was used
406  due to storage limitations in our computing environment. Then, one document was removed from
407  these sets of documents with its class noted. After that, the *positive* and *negative rules* under the
408  OCAT approach and the *positive* and *negative centroids* under the VSM were constructed using
409  the remaining 59 documents. In the third step the class of the document left out was inferred by
410  both algorithms. Then, the correctness of the classification was determined according to the cases
411  (a)–(h), as defined in Section 4. The previous second and third steps were repeated with different
412  sets of training examples until all 60 documents had their class inferred one at a time. This ex-
413  perimental setting was replicated 10 times with different subsets of the training data, at which
414  point the results of the two algorithms were tested for statistical differences.

415  *8.2. The 30/30 Cross-Validation*

416       The 30/30 Cross-Validation (30CV) was implemented on samples of 254 documents as follows.
417  The number of 254 documents was used to avoid excessive computational time. Initially, the
418  *positive* and *negative rules* under the OCAT approach and *positive* and *negative centroids* under the
419  VSM were constructed by using only 30 documents (randomly selected) from each class. Then, the
420  classification of the remaining 194 documents was inferred. As before, the correctness of this
421  classification was determined according to the cases (a)–(g), as defined earlier. As with the first
422  experimental setting, the 30CV validation was replicated 10 times, at which point the results of the
423  two algorithms were tested for statistical difference.

424  *8.3. Statistical performance of both algorithms*

425       To determine the statistical performance of both algorithms, the following hypotheses were
426  tested. The first test was needed to determine the relative dominance of the algorithms. In the
427  second test we implemented a sign test in order to determine the consistency of the dominance of
428  the algorithms.

14                  *S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*

1.  $H_0 : \bar{P}_{\text{OCAT}} \leqslant \bar{P}_{\text{VSM}}$

    $H_1 : \bar{P}_{\text{OCAT}} \leqslant \bar{P}_{\text{VSM}}$

2.  $H_0 : p = 0.50$

    $H_1 : p \neq 0.50$

431 where $P_{\text{OCAT}}$ and $P_{\text{VSM}}$ are the numbers of documents with "correct" classification under the two
432 algorithms divided by the total number of documents in the experiment. In addition, $p$ is the
433 probability of finding an equal number of positive and negative differences in a set of outcomes.
434 More on how these tests were performed is provided in the following sub-sections which present
435 the computational results.

436 *8.4. Experimental setting for the GLA*

437    Consider the question: *What is the best next document to be given to the oracle in order to*
438 *improve the performance of the classification rule*? Three samples of 510 documents (255 from each
439 class) from the three class-pairs: (DOE vs. ZIPFF), (AP. vs. DOE), and (WSJ vs. ZIPFF) were
440 used. The number of 510 documents was determined by the available RAM memory on the
441 Windows PC we used. The previous three class-pairs were processed by the OCAT algorithm
442 (only) under the RANDOM and the GUIDED learning approaches.
443    These two learning approaches were implemented as follows. At first, 30 documents from each
444 class in the experiment were randomly selected, and the positive and negative rules (logical ex-
445 pressions) were constructed. Next, the class membership of all 510 documents in the experiment
446 was inferred based on the two classification rules. The criteria expressed as cases (a)–(g) in Section
447 4 were used to determine the number of "correct", "incorrect", and "undecided" classifications.
448 Next, a document was added to the initial training sample as follows. For the case of the
449 RANDOM approach, this document was selected at random from among the documents not
450 included in the training sets yet (i.e., neither in $E^+$ nor in $E^-$).
451    In contrast, under the GUIDED approach this document was selected from the set of docu-
452 ments which the positive and negative rules had already termed as "undecided". However, if the
453 two rules did not detect an "undecided" case, then the GUIDED approach was replaced by the
454 RANDOM approach until a new "undecided" case was identified. This process for the RAN-
455 DOM and GUIDED approaches was repeated until all 510 documents were included in the two
456 training sets $E^+$ and $E^-$. The results of this experimentation are presented next.

457 **9. Results for the Leave-One-Out and the 30/30 Cross-Validations**

458    Table 3 summarizes the experimental results for the CV validation, while Table 4 summarizes
459 the results for the 30CV validation. The abbreviations "C:", "I:", and "U:" in the first column of
460 both tables correspond to the "correct", "incorrect", and "undecided" classification outcomes
461 which can be obtained by using the positive and the negative rules (for the OCAT case) or the
462 positive and the negative centroids (for the VSM case). For instance, the data in Table 3, column 2

Table 3
Summary of the first experimental setting: Leave-One-Out Cross-Validation

|  | DOE vs. AP | | DOE vs. WSJ | | DOE vs. ZIPFF | | AP vs. WSJ | | WSJ vs. ZIPPF | | TOTAL | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT |
| C: | 334 | 410 | 286 | 296 | 280 | 358 | 316 | 365 | 286 | 303 | 1502 | 1732 |
| I: | 261 | 5 | 314 | 66 | 320 | 25 | 284 | 47 | 314 | 76 | 1493 | 219 |
| U: | 5 | 185 | 0 | 238 | 0 | 217 | 0 | 188 | 0 | 221 | 5 | 1049 |

Table 4
Summary of the second experimental setting: 30/30 Cross-Validation

|  | DOE vs. AP | | DOE vs. WSJ | | DOE vs. ZIPFF | | AP vs. WSJ | | WSJ vs. ZIPPF | | TOTAL | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT | VSM | OCAT |
| C: | 975 | 1406 | 975 | 1266 | 1035 | 1320 | 1088 | 1283 | 1088 | 1145 | 5161 | 6420 |
| I: | 975 | 70 | 975 | 134 | 915 | 124 | 846 | 140 | 837 | 176 | 4548 | 644 |
| U: | 0 | 474 | 0 | 550 | 0 | 506 | 16 | 527 | 25 | 41 | 41 | 2686 |

463 (i.e., class-pair (DOE vs. AP)) indicate that the VSM identified 334 "correct", 261 "incorrect",
464 and 5 "undecided" cases.
465     Similarly, the data in Table 3, column 3 (i.e., class-pair (DOE vs. AP)) indicate that the OCAT
466 algorithm identified 410 "correct", 5 "incorrect", and 185 "undecided" classifications. The data in
467 the other columns can be interpreted in a similar manner. The last two columns of these two tables
468 summarize the results across all five class-pairs. Fig. 6 compares the proportions of these results
469 for both algorithms.
470     Two key observations can be derived from the size of the dark areas (or areas of "undecided"
471 classifications) in Fig. 6 which was derived from Tables 3 and 4. First, it can be observed that the
472 proportion of "undecided" cases detected by the VSM algorithm is almost 0%. These have oc-
473 curred when the two positive and the two negative centroids accepted the same document and,
474 therefore, the classes predicted by both sets of centroids have to be selected randomly. More
475 specifically, these "undecided" instances occurred even when these randomly predicted classes
476 were identical. The VSM was implemented using the CC coefficient, following the suggestions in
477 Salton (1989).
478     In contrast, as the second observation, we have large proportions of "undecided" classifications
479 with the OCAT algorithm. Please recall that this type of classification decision is desired because
480 such cases demonstrate that either the positive or the negative rules are unable to classify correctly
481 new documents. Therefore, in these results the large dark areas in the above figure show that both
482 rules were unable to classify correctly a large proportion of the documents in the experiments.
483 More importantly, the size of these areas indicates that positive or negative rules may be improved
484 if they are modified when an "undecided" situation is detected.
485     Consider the proportion of the "incorrect" classifications (i.e., the white areas in Fig. 6). One
486 can derive two conclusions. First, the number of "incorrect" classifications the VSM made
487 amounts to 49.77% (or 1493/3000) with the CV validation and to 46.65% (or 4548/9750) with the
488 30CV validation. These large proportions of "incorrect" classifications can be attributed to the
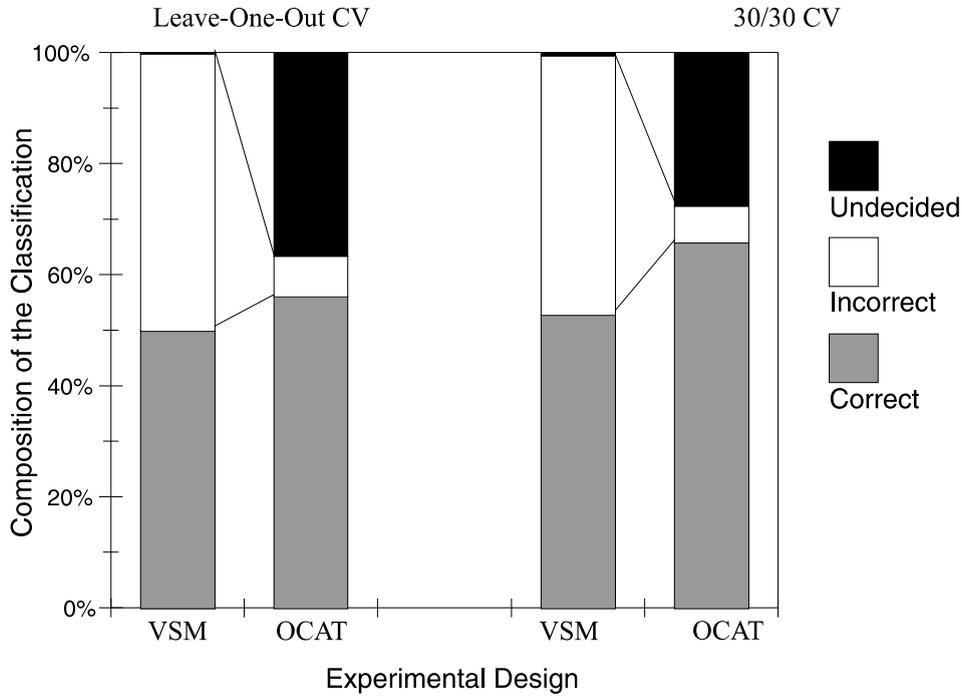
16                    *S. Nieto Sánchez et al. | Information Processing and Management xxx (2002) xxx–xxx*



Fig. 6. Comparison of the classification decisions under the OCAT and VSM approaches.

inability of the positive and negative centroids to distinguish between "incorrect" and "unde-cided" classifications. Second we have 7.30% (or 219/3000) and 6.61% (or 644/9750) of "incorrect" classifications the OCAT algorithm made with the two test settings. In this case, these small error rates can be attributed to the utilization of positive and negative rules of small size that enabled the OCAT algorithm to distinguish between the "incorrect" and "undecided" classifications.

Despite the disparate proportions of the "inaccurate" and "undecided" classifications for both of these algorithms, their performances were statistically compared using *only* the proportions with the "correct" classifications. That is, the undecided cases were not considered here. In this way the VSM approach was not placed in an unfair setting when it was compared with the OCAT approach. This comparison was needed in order to determine which algorithm better addressed

Table 5
Statistical difference in the classification accuracy of the VSM and OCAT approaches

| Type of experiment | $P_{OCAT}$[a] | $P_{VSM}$[b] | $P_{VSM} - P_{OCAT}$ | Binomial test | |
|---|---|---|---|---|---|
| | | | | Half-length[c] | Interval |
| CV | 0.577 | 0.501 | −0.0760 | 0.025 | (−0.035, −0.085) |
| 30CV | 0.658 | 0.529 | −0.1287 | 0.014 | (−11.47, −14.27) |

[a] 1732/$n$ and 6420/$n$; where $n$ is 3000 for CV and 9750 for 30CV.
[b] 1502/$n$ and 5161/$n$; where $n$ is 3000 for CV and 9750 for 30CV.
[c] Denotes that both approaches performed statistically differently.

Table 6
Data for sign test to determine the consistency in the ranking of the OCAT and VSM approaches

|  | Type of experiment | |
|---|---|---|
|  | CV | 30CV |
| Number of "+" signs | 4 | 7 |
| Number of "−" signs | 46 | 43 |
|  | $p\text{-value} = 2.23 \times 10^{-10}$ | $p\text{-value} = 1.04 \times 10^{-7}$ |

$$p\text{-value} = \sum_{i=0}^{m} \binom{50}{i} \cdot p^i \cdot (1-p)^{50-i},$$

where $m = 4$ for CV and 7 for 30CV and $p = 0.50$

499 the classification problem studied in this paper. For this comparison, it was assumed that no
500 additional improvement of the two algorithms was possible under the CV and 30CV cross-vali-
501 dations. Furthermore, the "incorrect" and "undecided" outcomes were considered as incorrect
502 classifications.
503   The results of these tests (as shown in Table 5) indicate that the OCAT approach is more
504 accurate in both types of computational experiments than the VSM. Furthermore, the very low *p*-
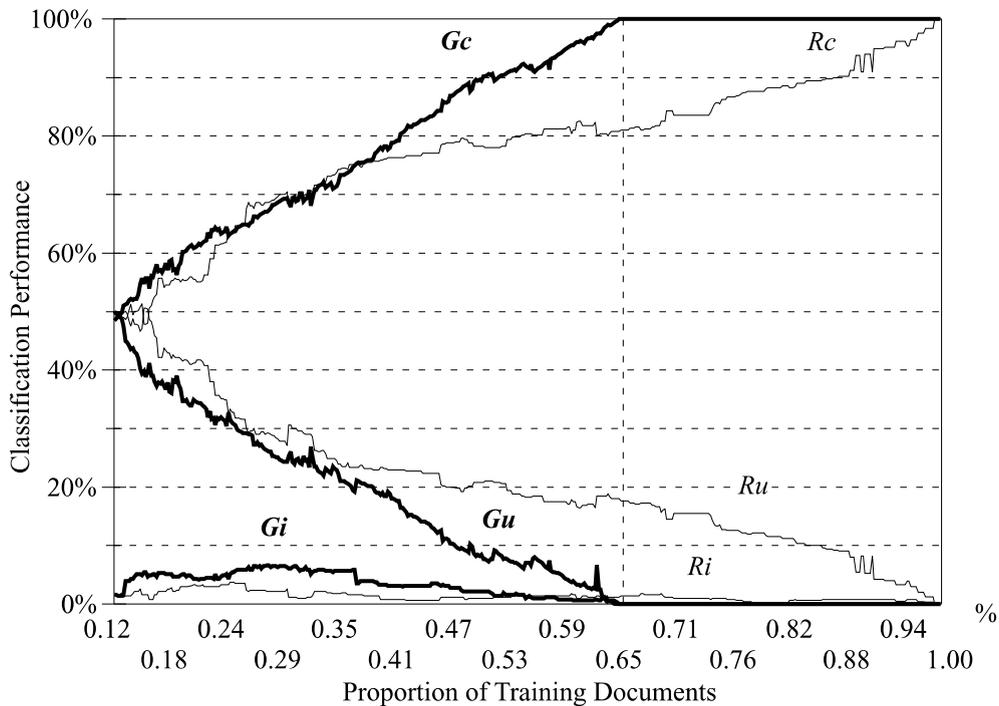


Fig. 7. Results when the GUIDED and RANDOM approaches were used on the (DOE vs. ZIPFF) class-pair.

505 values in Table 6 indicate that it is extremely unlikely to find a similar number of positive and
506 negative differences in the proportions of the "correct" classifications under the two approaches
507 (Barnes, 1994). Therefore, the results of these two statistical tests indicate that the OCAT ap-
508 proach is better suited to address the document classification problem studied in this paper.

509 **10. Results for the GLA**

510 Figs. 7–9 show the results of the OCAT algorithm under the RANDOM and GUIDED input
511 learning approaches. The horizontal axis indicates the percentage of training documents used
512 during the experiment. For example, at the beginning of the experiment there were 60 training
513 documents or 11.76% of the 510 documents in the experiment. Next, when one more document
514 was added to the training set, following the recommendation of the GUIDED and RANDOM
515 approaches, there were 12.16% of the documents in the experiment.
516 The vertical axis shows the proportions of "correct", "incorrect", and "undecided" classifi-
517 cation for the various percentages of training documents used in the experiment. The abbrevia-
518 tions Rc, Ri, Ru and Gc, Gi, Gu stand for the proportions of "correct", "incorrect", and
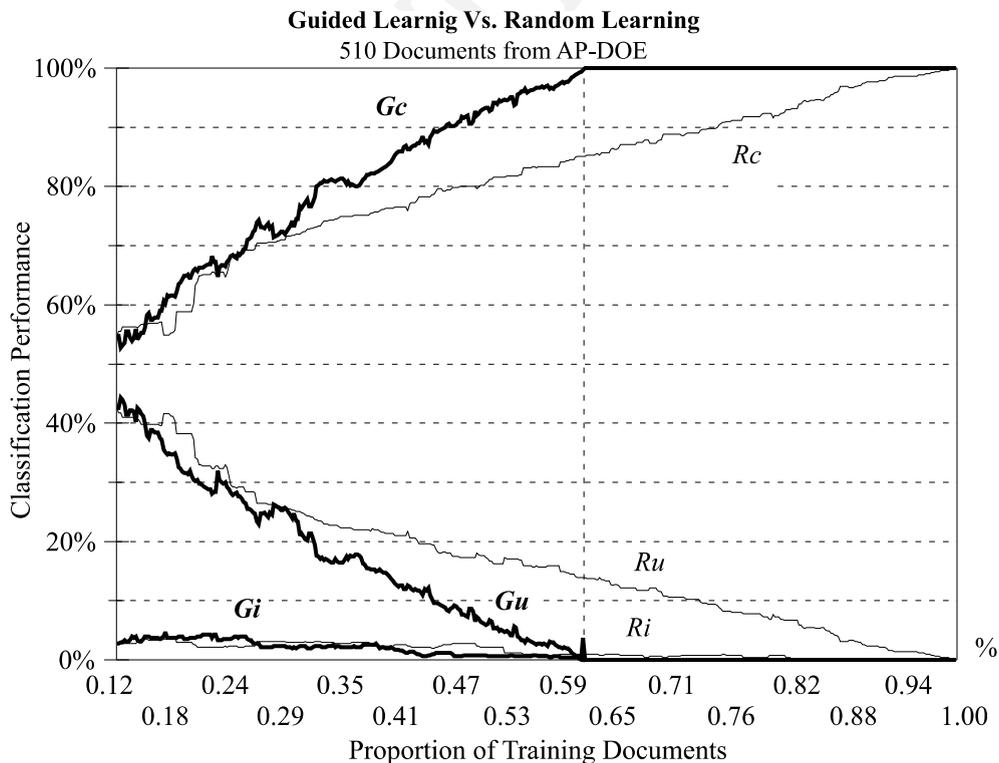519 "undecided" outcomes for the RANDOM and GUIDED approaches, respectively. For instance,



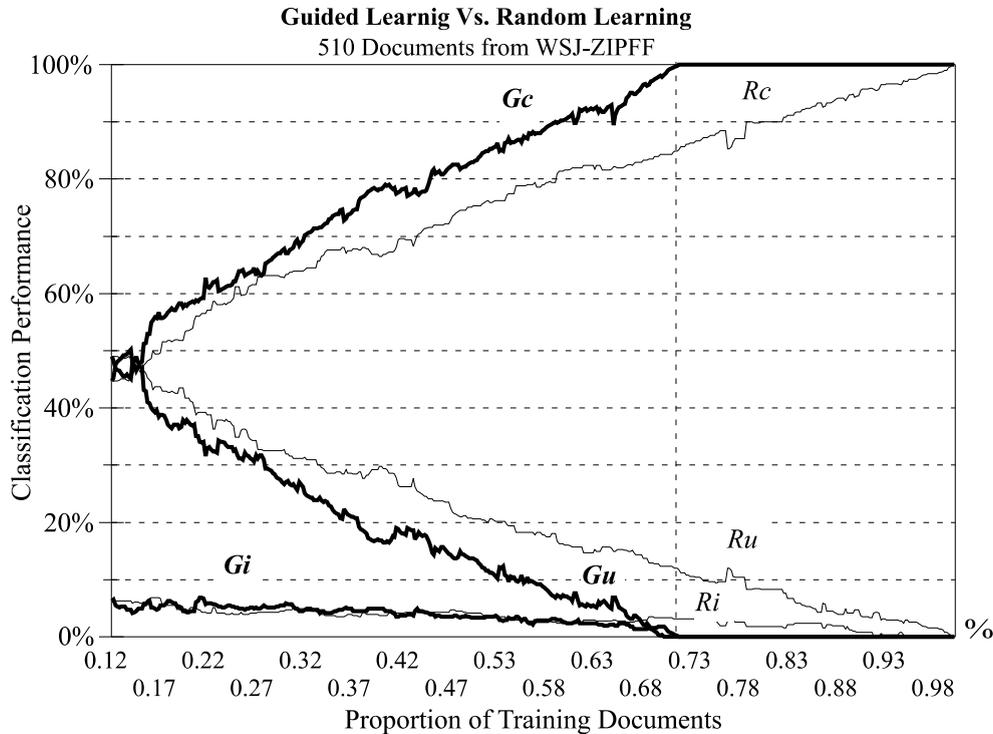Fig. 8.  Results when the GUIDED and RANDOM approaches were used on the (AP vs. DOE) class-pair.

**Guided Learnig Vs. Random Learning**
510 Documents from WSJ-ZIPFF



Fig. 9. Results when the GUIDED and RANDOM approaches were used on the (WSJ vs. ZIPFF) class-pair.

520 Rc is the proportion of "Correct" classifications under the RANDOM approach, and Gu is the
521 proportion of "Undecided" classification under the GUIDED approach.
522 Table 7 shows the percentage of training documents the OCAT algorithm needed before it
523 classified all 510 documents in each class-pair correctly (i.e., it became 100% accurate). The po-
524 sition of the dotted line in the above three figures corresponds to the percentages shown in this
525 table. For instance, in Fig. 7 (class-pair (DOE vs. ZIPFF)) this line is at 65.69% on the horizontal
526 axis.
527 Some important observations can be made regarding the proportions of "correct", "incorrect",
528 and "undecided" classifications in Figs. 7–9. First, the rate of "correct" classifications under the
529 GUIDED approach, Gc, was higher than the rate Rc under the RANDOM approach. Actually,

Table 7
Percentage of documents from the population that were inspected by the oracle before an accuracy of 100% was reached

| Class-pairs | % Under GUIDED | % Under RANDOM |
|---|---|---|
| (DOE vs. ZIPFF) | 65.69 | 100.00 |
| (AP vs. DOE) | 60.98 | 99.80 |
| (WSJ vs. ZIPFF) | 71.18 | 99.80 |
| Average | 65.95 | 99.87 |

100% accuracy was achieved when the number of "Incorrect" and "Undecided" classifications were 0%.

530  the last row in Table 7 indicates that the OCAT algorithm needed on the average about 34% less
531  training documents to classify correctly all 510 document under the GUIDED approach than
532  under the RANDOM approach.
533      These results are very interesting for a number of reasons. They confirm the assumption stated
534  in Section 5 which indicated that the utilization of documents with the "undecided" classification
535  could increase the accuracy of the OCAT algorithm. These results are also encouraging because
536  they help to answer the second question stated in the introduction of this section. That is, queries
537  to the oracle should stop when about 66% of the 510 documents from the three class-pairs of the
538  TIPSTER collection had been inspected and were included in the training sets. More importantly,
539  these results are important because they suggest that the OCAT algorithm can be employed for
540  the classification of large collections of text documents.
541      The other two observations are related to the rates at which the "incorrect" and "undecided"
542  classifications were eliminated. It can be observed from the previous three figures that these rates
543  were a direct consequence of improving the classification rules. The figures show that the rates Gi
544  and Gu reach 0% when about 66% (336 documents) of the 510 documents in the experiment have
545  been included in $E^+$ and $E^-$. On the other hand, it can be seen that under the RANDOM learning
546  approach, the rates Ri and Ru reached 0% when 99.8% (509 documents) of the documents are
547  processed.

548  **11. Concluding remarks**

549      This paper has examined a classification problem in which a document must be classified into
550  one of two disjoint classes. As an example of the importance of this type of classification, one can
551  consider the possible release to the public of documents that may affect national security. The
552  method proposed in this paper (being an automatic method) is not infallible. This is also true
553  because its performance depends on how representative the training examples (documents) are.
554  The application of such an approach to a problem of critical importance (such as the one high-
555  lighted in the introduction) can be seen as an important and useful automatic tool for a pre-
556  liminary selection of the documents to be classified.
557      We considered an approach to this problem based on the VSM algorithm and compared it with
558  an algorithm which is based on mathematical logic, called the OCAT algorithm. We tested these
559  two approaches on almost 3000 documents from the four document classes of the TIPSTER
560  collection: Department of Energy (DOE), Wall Street Journal (WSJ), Associated Press (AP), and
561  the ZIPFF class. Furthermore, these documents were analyzed under two types of experimental
562  settings: (i) Leave-One-Out Cross-Validation and (ii) a 30/30 Cross-Validation (where 30 indicates
563  the initial number of training documents from each document class). The experimental results
564  suggest that the OCAT algorithm performed significantly better in classifying documents into two
565  disjoint classes than the VSM.
566      Moreover, the results of a third experiment suggested that the classification efficiency of the
567  OCAT algorithm can be improved substantially if a GLA is implemented. Actually, experiments
568  on samples of 510 documents from the previous four classes of the TIPSTER collection indicated
569  that the OCAT algorithm needed only about 336 (i.e., 66% of the) training documents before it
570  correctly classified all of the documents.

IPM 564
DISK / 28/1/02

ARTICLE IN PRESS

No. of pages: 22
DTD 4.3.1/ SPS-N

*S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*          21

571 The results presented here, although limited to a relatively small collection of almost 3000
572 documents, are encouraging because they suggest that the OCAT algorithm can be used in the
573 classification of larger collections of documents.

## 574 Acknowledgements

## 583 References

584 Aldenderfer, M. S. (1984). *Cluster analysis*. Beverly-Hill, CA: Sage.
585 Anderberg, M. R. (1973). *Cluster analysis for applications*. New York: Academic Publishers.
586 Barnes, J. W. (1994). *Statistical analysis for engineers and scientists, a computer-based approach*. New York: McGraw-
587     Hill.
588 Buckley, C., & Salton, G. (1995). Optimization of relevance feedback weights. In *Proceedings of SIGIR 1995* (pp. 351–
589     357).
590 Chen, H. (1996). *Machine learning approach to document retrieval: An overview and an experiment*. Technical Report,
591     University of Arizona, MIS Department, Tucson, AZ, USA.
592 Cleveland, D., & Cleveland, A. D. (1983). *Introduction to indexing and abstracting*. Littleton, CO: Libraries Unlimited.
593 Deshpande, A. S., & Triantaphyllou, E. (1998). A greedy randomized adaptive search procedure (GRASP) for inferring
594     logical clauses from examples in polynomial time and some extensions. *Mathematical and Computer Modelling*,
595     *27*(1), 75–99.
596 DOE (1995). *General Course on Classification/Declassification, Student Syllabus, Handouts, and Practical Exercises*. US
597     Department of Energy, Germantown, MD, USA.
598 DynMeridian (1996). *Declassification Productivity Initiative Study Report*. DynCorp Company, Report Prepared for the
599     US Department of Energy, Germantown, MD, USA.
600 Fox, C. (1990). A stop list for general text. *ACM Special Interest Group on Information Retrieval*, *24*(1–2), 19–35.
601 Harman, D. (1995). Overview of the second text retrieval conference (TREC-2). *Information Processing and
602     Management*, *31*(3), 271–289.
603 Luhn, H. P. (1958). The automatic creation of literature abstracts. *IBM Journal of Research and Development*, *5*(3),
604     155–165.
605 Luhn, H. P. (1957). A statistical approach to mechanized encoding and searching of literary information. *IBM Journal
606     of Research and Development*, *4*(4), 600–605.
607 Meadow, C. T. (1992). *Text information retrieval systems*. San Diego, CA: Academic Press.
608 Nieto Sanchez, S., Triantaphyllou, E., Chen, J., & Liao, T. W. (2001). An incremental learning algorithm for
609     constructing Boolean functions from positive and negative examples. *Computers and Operations Research* (in press).
610 Salton, G. (1968). *Automatic information organization and retrieval*. New York: McGraw-Hill.
611 Salton, G., & Wong, A. (1975). A vector space model for automatic indexing. *Information retrieval and language
612     processing*, *18*(11), 613–620.

IPM 564
DISK / 28/1/02

ARTICLE IN PRESS

No. of pages: 22
DTD 4.3.1/ SPS-N

22                    *S. Nieto Sánchez et al. / Information Processing and Management xxx (2002) xxx–xxx*

613  Salton, G. (1989). *Automatic text processing. The transformation analysis, and retrieval of information by computer*.
614      Reading, MA: Addison-Wesley.
615  Scholtes, J. C. (1993). *Neural networks in natural language processing and information retrieval*. The Netherlands: North-
616      Holland.
617  Späth, H. (1985). *Cluster dissection and analysis: theory, Fortran programs, and examples*. Chichester, UK: Ellis
618      Harwood.
619  Shaw, W. M. (1995). Term-relevance computations and perfect retrieval performance. *Information Processing and*
620      *Management*, *31*(4), 312–321.
621  Triantaphyllou, E. (2001). *The OCAT approach for data mining and knowledge discovery*, Working Paper, IMSE
622      Department, Louisiana State University, Baton Rouge, LA 70803-6409, USA.
623  Triantaphyllou, E., & Soyster, A. L. (1996a). On the minimum number of logical clauses inferred from examples.
624      *Computers and Operations Research*, *23*(8), 783–799.
625  Triantaphyllou, E., & Soyster, A. L. (1996b). An approach to guided learning of Boolean functions. *Mathematical*
626      *Computing Modeling*, *23*(3), 69–86.
627  Triantaphyllou, E., & Soyster, A. L. (1995). A relationship between CNF and DNF systems which are derived from the
628      same positive and negative examples. *ORSA Journal on Computing*, *7*(3), 283–285.
629  Triantaphyllou, E., Soyster, A. L., & Kumara, S. R. T. (1994). Generating logical expressions from positive and
630      negative examples via a branch-and-bound approach. *Computers and Operations Research*, *21*(2), 185–197.
631  Triantaphyllou, E. (1994). Inference of a minimum size Boolean function from examples by using a new efficient
632      branch-and-bound approach. *Journal of Global Optimization*, *5*, 64–94.
633  Van Rijsbergen, C. J. (1979). *Information retrieval* (2nd ed.). London: Butterworths.
634  Voorhees, E. (1998). Overview of the sixth text retrieval conference (TREC-6). In *Proceedings of the sixth text retrieval*
635      *conference (TREC-6), Gaithersburg, MD, USA* (pp. 1–27).
636  Zipff, H. P. (1949). *Human behavior and the principle of least effort*. Menlo Park, CA: Addison-Wesley.