

HIERARCHICAL ENSEMBLE LEARNING FOR MULTIMEDIA CATEGORIZATION AND AUTOANNOTATION

Serhiy Koisnov and Stéphane Marchand-Maillet
Computer Vision and Multimedia Lab, University of Geneva
24 rue du General-Dufour, Geneva, Switzerland
{kosinov,marchand}@cui.unige.ch

Abstract. This paper presents a hierarchical ensemble learning method applied in the context of multimedia autoannotation. In contrast to the standard multiple-category classification setting that assumes independent, non-overlapping and exhaustive set of categories, the proposed approach models explicitly the hierarchical relationships among target classes and estimates their relevance to a query as a trade-off between the goodness of fit to a given category description and its inherent uncertainty. The promising results of the empirical evaluation confirm the viability of the proposed approach, validated in comparison to several techniques of ensemble learning, as well as with different type of baseline classifiers.

INTRODUCTION

One of the essential challenges in modern information retrieval is to be able to deduce high-level semantics from the low-level perceptual features of multimedia, which the literature sources refer to as the semantic categorization, keyword prediction, autoannotation or automatic linguistic indexing task. The diversity in the problem terminology reflects the variety of contributions from numerous research domains that have been proposed to date. For example, an appealing idea of treating the visual feature data as another language to translate semantic keywords to and from is developed with the aid of generative probabilistic models by Barnard *et al.* [1, 2]. A family of methods [16, 18, 22, 23], related to the cross-language extension of the latent semantic indexing (LSI) technique [5, 11], permit the retrieval of multimedia semantics via low-level feature queries. Yet, the majority of the other approaches consider the multimedia autoannotation problem in the multiple-category classification framework, where unseen documents must be assigned to one or more predefined semantic categories. In [7], for instance, the authors focus on improving several popular ensemble schemes, such as OPC (one per

class), PWC (pair-wise coupling) and ECOC (error-correcting output codes). The methods developed in [3, 12, 13] decompose a multiple-category classification task into a collection of binary classification problems and propose ways of recombining effectively the individual predictions from classifiers as diverse as SVM, BPM, 2D-MHMM. The semantic categories for these and many other classification-based techniques are generally assumed to be independent, non-overlapping and sufficient to cover all of the problem domain.

The approach presented in this paper is also formulated as a classification-based method, but differs from the above work in the important respect that the relationships among the semantic categories derived from the individual keywords of the annotation corpora are explicitly modeled in Bayesian terms, leading to a more consistent autoannotation performance. Furthermore, the proposed method broadens the range of the derived annotation allowing to predict more general notions or semantically-related keyword groups in addition to individual keywords present in the training data vocabulary. Another benefit of the proposed formulation is that it gives an answer to such an important question as how many keywords the system should predict and whether it is reasonable to predict anything at all.

The remainder of this paper is organized as follows. Section 2 presents the problem formulation focused on the autoannotation of digital images as a particular form of multimedia documents, followed by an illustrative example of the proposed method, given in Section 3. The experimental results and concluding remarks are provided in Sections 4 and 5.

PROBLEM FORMULATION

We employ a hierarchical ensemble of binary classifiers in order to perform semantic annotation of unseen images. Given a training set of annotated images $X^{(T)} = \{I_t, K_t\}_{t=1}^n$, where I_t and K_t represent the feature vector of a given image and its associated set of keywords, respectively, the concept hierarchy $H = \{C_i\}_{i=1}^N$ is defined by all of the unique nouns comprising the annotation vocabulary $V = \bigcup_{t=1}^n K_t$ and their hyponyms derived from WordNet [15]. Every concept C_i occupies a separate node in H , and is associated with a binary classifier Φ_i designed to distinguish the set of leaf concepts subsumed (directly or indirectly) by C_i , denoted as $\mathbf{L}(C_i)$, from all of the others. An example of a hierarchy derived for a simple vocabulary $V: \{beach, flower, grass, mountain, rock, sky, tree\}$ is shown in Figure 1.

In order to perform the autoannotation of an unseen image represented by a low-level feature vector I_U , each concept C_i is assessed as a potential candidate. Thus, the set of possible annotations is no longer restricted to be V , as is the case for the majority of other similar techniques. The relevance of C_i is seen as a trade-off between, on one hand, how well the input data I_U fits the description of C_i from the classification accuracy point of view, and, on the other hand, how specific or non-ambiguous the candidate set of keywords $\mathbf{L}(C_i)$ is. In our method, the first of these two quantities is represented by

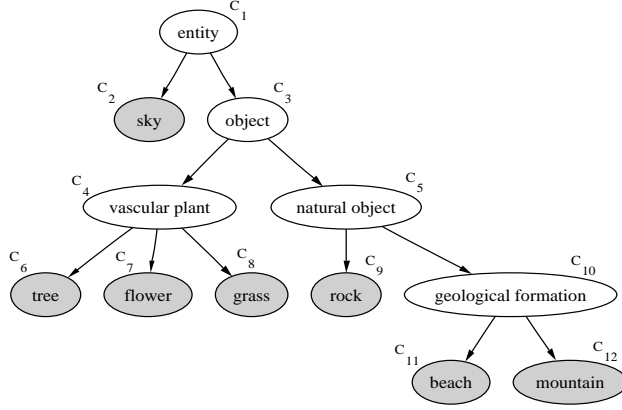


Figure 1: Classifier hierarchy example. Shaded nodes denote $C_i \in V$

the posterior probability of a concept given the data, $P(C_i|I_U)$, while the second one is estimated as the posterior probability of a concept given the assumption that a particular keyword k from the set of all homonyms of C_i is chosen correctly, denoted as $P(C_i|k)$.

For a given concept C_i , the estimate of $P(C_i|I_U)$ is determined according to the following theorem, which is a reformulation of a previously established result described in [10]:

Theorem 1, (Kumar et al., 2002). *The posterior probability $P(C_i|I_U)$ for any input I_U is the product of the posterior probabilities of all the internal classifiers along a unique path from the root node to C_i , i.e.*

$$P(C_i|I_U) = \prod_{l=0}^{\mathcal{D}(C_i)-1} P(C_i^{(l+1)}|I_U, C_i^{(l)}), \quad (1)$$

where $\mathcal{D}(C_i)$ is the depth of C_i (the depth of the root concept C_1 is 0), $C_i^{(l)}$ is the concept at depth l on the path from the root node to C_i , such that $C_i^{(\mathcal{D}(C_i))} \equiv C_i$ and $C_i^{(0)} \equiv C_1$.

In order to ensure that (1) is applicable in the case of classifiers with non-probabilistic outputs, such as SVM [4], a sigmoid function, e.g., $\frac{1}{1+\exp(Ay_i+B)}$, is fit to the raw classifier output values y_i , as described in [17]. As for $P(C_i|k)$, the Bayes theorem allows to express this quantity in terms of statistics of the training data as shown in (2):

$$P(C_i|k) = \frac{P(k|C_i)P(C_i)}{\sum_{C_i \in H} P(k|C_i)P(C_i)}, \quad (2)$$

where $P(C_i)$, a prior probability of concept C_i , is estimated from the training data as:

$$P(C_i) = \frac{\sum_{C \in \mathbf{L}(C_i)} \text{freq}^{(T)}(C)}{\sum_{C \in V} \text{freq}^{(T)}(C)}, \quad (3)$$

and $P(k|C_i)$, the worst-case estimate of the probability of choosing a correct annotation keyword k given the degree of generality of concept C_i , is deduced from the homonym set cardinality information derived from WordNet:

$$P(k|C_i) = \frac{\min_{C \in \mathbf{L}(C_i)} \text{freq}^{(W)}(C)}{\text{freq}^{(W)}(C_i)}. \quad (4)$$

In (3) and (4), the frequency of a given concept in the training data and the cardinality of the WordNet homonym set are denoted as $\text{freq}^{(T)}$ and $\text{freq}^{(W)}$, respectively.

Finally, assuming that the likelihood of the input data I_U given C_i is not dependent on the correctness of a particular choice of k from the homonym set of C_i , we obtain the following result:

$$P(C_i|I_U, k) \propto P(C_i|I_U)P(C_i|k) = \rho, \quad (5)$$

which essentially represents a means of comparison of different hypothesis concepts $\{C_i\}$ that takes into account both the goodness of fit of the data I_U to a given concept description and the concept's inherent degree of uncertainty or specificity. The next section illustrates these notions.

ILLUSTRATIVE EXAMPLE

Let us come back to the simplified 12-concept classifier hierarchy given in Figure 1. To be able to observe the effect of each of the two factors contributing to the final estimate of the concept relevance, ρ , we plot separately the computed values of $P(C_i|k)$, Figure 2(c), and $P(C_i|I_U)$, Figure 2(b), for a sample test image query depicted in Figure 2(a). As the diagrams show,

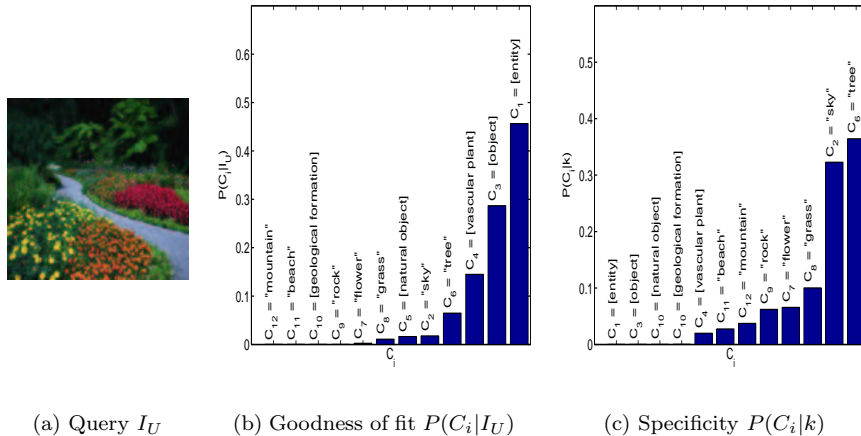


Figure 2: Individual contributions of factors $P(C_i|I_U)$ and $P(C_i|k)$

there is a natural tendency among the values of $P(C_i|I_U)$ to favor simpler, more general concepts, such as *object*, due to the smaller number of terms to be evaluated in product (1). Quite the opposite trend is noticeable among the estimates of $P(C_i|k)$ that tend to promote very specific, unambiguous concepts, such as *sky*, taking into account their prior probabilities as well. This very trade-off of “Goodness of fit vs. Specificity” is captured by the concept relevance, ρ , leading to the results listed in Table 1 that demonstrate a

TABLE 1: CANDIDATE CONCEPTS RANKED BY RELEVANCE

Rank	$-\log_2 \rho(C_i)$	Concept C_i
1	5.41	$C_6 = \text{tree}$
2	7.46	$C_2 = \text{sky}$
3	8.44	$C_4 = \text{vascular plant (flower, grass, tree)}$
4	9.84	$C_8 = \text{grass}$
5	12.64	$C_7 = \text{flower}$
6	17.26	$C_1 = \text{entity}$
7	17.87	$C_3 = \text{object}$
8	19.42	$C_5 = \text{natural object}$
9	21.00	$C_9 = \text{rock}$
10	44.32	$C_{10} = \text{geological formation}$
11	55.97	$C_{12} = \text{mountain}$
12	56.35	$C_{11} = \text{beach}$

reasonable degree of coherence between the top ranking concepts C_i and the true keywords of the query $K_U = \{\textit{flowers, path, grass, trees}\}$.

Another important property of the proposed method that the figures from Table 1 help highlight is its ability to determine exactly how many of the top-ranked concepts should be predicted. Many existing approaches [1, 2, 16] resolve this issue by specifying a tunable “refuse-to-predict” parameter that regulates the propensity of image regions to emit concepts or, as some other techniques, by simply considering a fixed number of top-ranked entries. In our case, the relevance of the root node, $\rho_1 = \rho(C_1)$, provides a natural threshold that determines the number of candidate annotation concepts to be selected. An intuitive interpretation of neg-logarithm of this quantity comes from the minimum message length (MML) principle of information theory [21], which interprets $-\log_2 \rho_1$ as the null-model hypothesis test that corresponds to transmitting all the data, since the root concept subsumes all of the other concepts, as is. According to the MML principle, any hypothesis that cannot better the null-model is not acceptable. In our example, this assertion makes us discard all of the candidate concepts ranked 6 or worse (see Table 1).

EXPERIMENTAL RESULTS

In our experiments we have used data from two separate image collections for training and testing in an attempt to ensure collection-independent learning. The training data was derived from the Washington University annotated image collection [14] with about 600 images, while the testing data constituted a 254 image subset, New Zealand and Ireland sections, from Corel image database. The visual information for each training image was represented by 286-dimensional feature vector containing 166 global color histogram and 120 Gabor filter texture descriptors extracted by the *Viper* system [19]. Annotation keywords appearing only once were eliminated from the target vocabulary V , from which a hierarchical ensemble of 60 concepts was constructed.

In order to be able to judge the performance of the presented method in terms of the traditional precision and recall indicators, we have adopted the following strategy. Whenever a non-leaf concept, $C_i \notin V$, is predicted, it is evaluated as a union of its underlying keywords, $\mathbf{L}(C_i)$, thus bridging the vocabulary gap between the derived concepts, e.g. [*vessel*, *watercraft*], and the actual training data, e.g. *boat*, *sailboat*, *ferryboat*, *rowboat*, at the expense of precision. Using the DDA baseline classifiers [8, 9] for each concept $C_i \in H$, the following precision and recall results on the test set vocabulary were obtained (see Figure 3). As seen from the figure, the naturally high recall

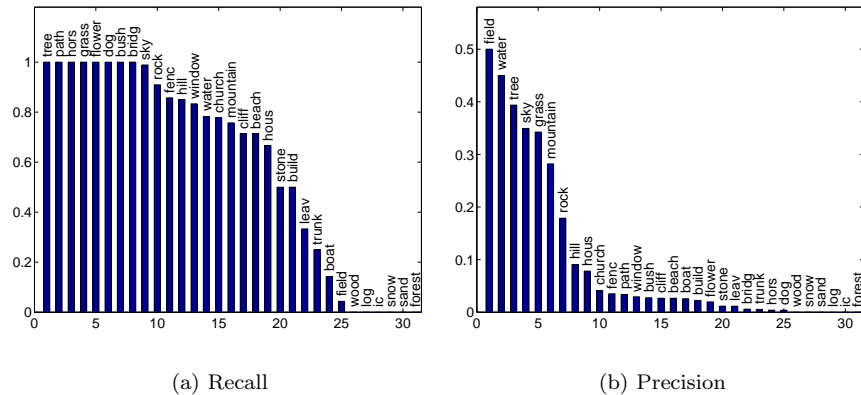


Figure 3: Performance indicators on test data vocabulary

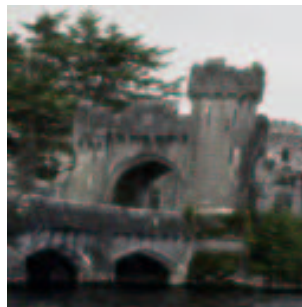
results boosted by keyword group retrieval, Figure 3(a) do not necessarily correspond to high frequency common concepts emphasizing the importance of the concept co-occurrence factors, while the significantly lower precision values for complex concepts, such as *church*, *fence*, *boat*, Figure 3(b), indicate that these words are much more often retrieved as a group of semantically-related keywords, rather than individually.

An illustration of the automatically derived annotation is provided in Figure 4, showing examples occurrences of out-of-vocabulary words being

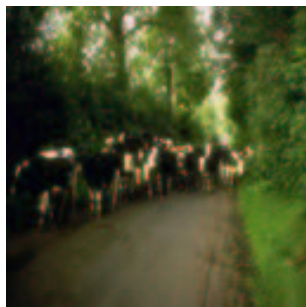
replaced by a visually similar common concepts $C_i \in V$ (top-right image, *castle* \rightarrow *rock*), members of the vocabulary being predicted as semantically relevant, but more common (and therefore, more likely) concepts C_i (top-left, *buildings* \rightarrow *construction*), as well as other typical predictions.



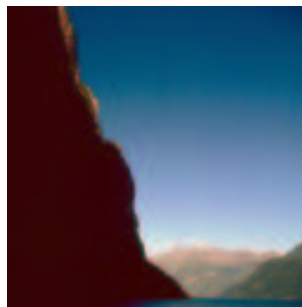
True annotation:
sky, street, buildings, town
Autoannotation:
sky, construction,
natural object, artefact



True annotation:
sky, castle, water, tree
Autoannotation:
sky, rock, tree



True annotation:
cows, road, trees, grass
Autoannotation:
bush, tree, grass, vascular
plant, woody plant, organism



True annotation:
sky, water, mountain, trees,
Autoannotation:
sky, water, geological formation,
natural object, artefact

Figure 4: Autoannotation of test images

In addition to the above experiments, we have compared the presented method to several popular classifier ensemble techniques, such as OPC, or one-against-all strategy, and Max Wins algorithms [6] that combined SVM baseline classifiers. As shown in Table 2, the proposed hierarchical semantic ensemble (HSE) approach achieved better results despite the fact that only a fixed number of top-ranked singleton concepts was allowed to be predicted, which was done in order to ensure equal conditions for all of the methods, most of which have no means of determining exactly the number of concepts in the derived annotation. The first row of Table 2 represents the reference point performance attained by sampling concepts according to their empirical distribution in the training data annotation, i.e. picking word *tree* first, since it is most likely to occur, then *sky*, and so on, whereas the last row shows

TABLE 2: CLASSIFIER ENSEMBLE PERFORMANCE RESTRICTED TO TOP 5 KEYWORDS

Ensemble	Baseline classifier	% Recall	% Precision
Empirical	none	16.13	5.04
Max Wins	SVM, polynomial	8.14	3.83
Max Wins	SVM, gaussian	10.61	4.47
OPC	SVM, polynomial	20.31	7.85
OPC	SVM, gaussian	21.27	10.19
HSE	DDA	21.22	10.20
HSE+siblings	DDA	28.42	26.88

an improvement in performance of the presented HSE method when one considers sibling concepts¹ the same, e.g. *sailboat* and *boat*.

We also examined the performance of various types of binary SVM techniques as baseline classifiers in the proposed HSE framework, as illustrated in Table 3. The results of these studies have confirmed earlier findings [20]

TABLE 3: INFLUENCE OF BASELINE CLASSIFIER ON HSE PERFORMANCE

Baseline classifier	% Recall	% Precision
SVM, linear	18.12	5.28
SVM, polynomial	18.34	5.67
SVM, gaussian	18.62	6.05
DDA	21.22	10.20

stating that state-of-the-art individual classifiers do not necessarily always lead to a better performance in ensembles, while the inadequate results for the Max Wins technique, the only scheme to be using raw classifier outputs, emphasize the importance of the role of fitted posterior probabilities in classification ensembles.

CONCLUSION

We have presented a hierarchical ensemble learning method applied in the context of multimedia autoannotation. In contrast to the standard multiple-category classification setting that assumes independent, non-overlapping and exhaustive set of categories, the proposed approach models explicitly the hierarchical relationships among target classes using WordNet, and estimates their relevance to a query as a trade-off between the goodness of fit to a given category description and its inherent uncertainty. The latter aspect, formulated in Bayesian terms, brings an additional benefit of allowing to determine exactly the number of categories to be predicted. The promising results of the empirical evaluation confirm the viability of the proposed approach, validated in comparison to several techniques of ensemble learning, as well as with different type of baseline classifiers.

¹Concept A is a sibling of concept B if $A^{(\mathcal{D}(A)-1)} = B^{(\mathcal{D}(B)-1)}$.

In perspective, we plan to explore the problem of establishing correspondence between individual annotation keywords and low-level feature descriptors, and improve the proposed approach by taking advantage of the meaningful structure of the resulting hierarchical classification ensemble in order to incorporate relevance feedback from the user.

ACKNOWLEDGEMENTS

This work is supported by Swiss National Fund for Scientific Research (SNF) under grant number 2100.066648 and the Swiss Interactive Multimodal Information Management (NCCR (IM)2) network.

REFERENCES

- [1] K. Barnard, P. Duygulu, N. de Freitas, D. Forsyth, D. Blei and M. Jordan, "Matching Words and Pictures," **Journal of Machine Learning Research**, vol. 3, pp. 1107–1135, 2003.
- [2] K. Barnard, P. Duygulu and D. Forsyth, "Recognition as Translating Images into Text," **Internet Imaging IX, Electronic Imaging 2003 (Invited paper)**, 2003.
- [3] E. Chang, K. Goh, G. Sychay and G. Wu, "CBSA: content-based soft annotation for multimodal image retrieval using Bayes point machines," in **IEEE Transactions on Circuits and Systems for Video Technology**, vol. 13, pp. 26–38, 2003.
- [4] N. Cristianini and J. Shawe-Taylor, **An introduction to Support Vector Machines and other kernel-based learning methods**, Cambridge University Press, 2000.
- [5] S. Deerwester, S. Dumais, G. Furnas, T. Landauer and R. Harshman, "Indexing by Latent Semantic Analysis," **Journal of the American Society of Information Science**, , no. 41, pp. 391–407, 1990.
- [6] J. Friedman, "Another approach to polychotomous classification," Techn. report, **Stanford University**, 1996.
- [7] K.-S. Goh, E. Chang and K.-T. Cheng, "SVM binary classifier ensembles for image classification," in **Proceedings of the tenth international conference on information and knowledge management**, ACM Press, 2001, pp. 395–402.
- [8] S. Koisnov, S. Marchand-Maillet and T. Pun, "Iterative majorization approach to the distance-based discriminant analysis," Presented by S. Koisnov at "Conference of the GfKI 2004", Dortmund, Germany.
- [9] S. Koisnov, S. Marchand-Maillet and T. Pun, "Visual object categorization using distance-based discriminant analysis," in **Proceedings of the 4th International Workshop on Multimedia Data and Document Engineering**, July 2004, to appear.
- [10] S. Kumar, J. Ghosh and M. M. Crawford, "Hierarchical Fusion of Multiple Classifiers for Hyperspectral Data Analysis," **Pattern Analysis and Applications**, vol. 5, pp. 210–220, 2002.

- [11] T. Landauer and M. Littman, "Fully automatic cross-language document retrieval using latent semantic indexing," in **Sixth Annual Conference of the UW Centre for the New Oxford English Dictionary and Text Research**, UW Centre for the New OED and Text Research, Waterloo, Ontario, 1990, pp. 31–38.
- [12] B. Li and K. Goh, "Confidence-based dynamic ensemble for image annotation and semantics discovery," in **Proceedings of the eleventh ACM international conference on Multimedia**, ACM Press, 2003, pp. 195–206.
- [13] J. Li and J. Z. Wang, "Automatic Linguistic Indexing of Pictures by a Statistical Modeling Approach," **IEEE Trans. Pattern Anal. Mach. Intell.**, vol. 25, no. 9, pp. 1075–1088, 2003.
- [14] Y. Li and L. G. Shapiro, "Object Recognition for Content-Based Image Retrieval," in **Lecture Notes in Computer Science**, Springer-Verlag, to appear, 2004.
- [15] G. A. Miller, "WordNet: a lexical database for English," **Commun. ACM**, vol. 38, no. 11, pp. 39–41, 1995.
- [16] F. Monay and D. Gatica-Perez, "On Image Auto-Annotation with Latent Space Models," in **Proc. ACM Int. Conf. on Multimedia (ACM MM)**, November 2003.
- [17] J. Platt, "Probabilistic Outputs for support vector machines and comparison to regularized likelihood methods," in A. Smola, P. Bartlett, B. Schölkopf and D. Schuurmans (eds.), **Advances in Large Margin Classifiers**, MIT Press, 1999.
- [18] P. Praks, J. Dvorsky and V. Snasel, "Latent Semantic Indexing for Image Retrieval Systems," in **Proceedings of the SIAM Conference on Applied Linear Algebra (LA03)**, The College of William and Mary, Williamsburg, USA, 2003.
- [19] D. M. Squire, W. Müller, H. Müller and J. Raki, "Content-based query of image databases, inspirations from text retrieval: inverted files, frequency-based weights and relevance feedback," in **The 11th Scandinavian Conference on Image Analysis**, Kangerlussuaq, Greenland, june 1999, pp. 143–149.
- [20] V. Tresp, "A Bayesian Committee Machine," **Neural Computation**, vol. 12, no. 11, pp. 2719–2741, 2000.
- [21] C. Wallace and D. Dowe, "Minimum Message Length and Kolmogorov complexity," **Computer Journal**, vol. 42, no. 4, pp. 270–283, 1999.
- [22] R. Zhao and W. Grosky, "From Features to Semantics: Some Preliminary Results," in **IEEE International Conference on Multimedia and Expo (II)**, 2000, pp. 679–682.
- [23] R. Zhao and W. Grosky, "Narrowing the semantic gap - Improved text-based document web document retrieval using visual features," **IEEE Trans. on Multimedia**, vol. 4, no. 2, pp. 189–200, 2002.