

Theoretical and Experimental Analysis of a Two-Stage System for Classification

Nicola Giusti, Francesco Masulli, *Member, IEEE*, and
Alessandro Sperduti, *Member, IEEE Computer Society*

Abstract—We consider a popular approach to multicategory classification tasks: a two-stage system based on a first (global) classifier with rejection followed by a (local) nearest-neighbor classifier. Patterns which are not rejected by the first classifier are classified according to its output. Rejected patterns are passed to the nearest-neighbor classifier together with the $\text{top-}h$ ranking classes returned by the first classifier. The nearest-neighbor classifier, looking at patterns in the $\text{top-}h$ classes, classifies the rejected pattern. An editing strategy for the nearest-neighbor reference database, controlled by the first classifier, is also considered. We analyze this system, showing that even if the first level and nearest-neighbor classifiers are not optimal in a Bayes sense, the system as a whole may be optimal. Moreover, we formally relate the response time of the system to the rejection rate of the first classifier and to the other system parameters. The error-response time trade-off is also discussed. Finally, we experimentally study two instances of the system applied to the recognition of handwritten digits. In one system, the first classifier is a fuzzy basis functions network, while in the second system it is a feed-forward neural network. Classification results as well as response times for different settings of the system parameters are reported for both systems.

Index Terms—Multicategory classification, rejection, global and local classification, hierarchical classifier, Bayes classifier.



1 INTRODUCTION

IN complex multicategory classification tasks it is widely used the approach where a multistage or hierarchical system is used in order to find the right trade-off between accuracy and resources allocation (e.g., response time, size of system, error costs). A typical example is a system where the set of classes is organized in a hierarchical way and different classifiers are trained in order to drive the input pattern towards the most specific classifier to be used for the final classification (see, for example, [22], [31], [9], [24], [30]). Another example is a system where only a subset of the input features, i.e., the less expensive to compute are given as input to a first-level classifier and further input features are eventually used at further classification stages if the final classification cannot be performed with a sufficient level of confidence (e.g., reject-based approaches using an incremental set of input features [14], [23], [13], [31], [25]).

In this paper, we focus on a scheme where a fast-first classifier with rejection is used to classify patterns with high confidence. Rejected patterns are forwarded to a more complex and slower second-level classifier for a final classification (or further rejection). Typically, the system as a whole holds better classification performance with respect to the first classifier at the cost of a slower response time. Alternatively, improved classification performance can also be obtained by resorting to a committee of

classifiers [5], [28], [10], [20], [26], [32], [21], [29], [7], [17], [16], [11], [27]; however, the above hierarchical system turns out to be more flexible if constraints on the mean response time are imposed by the operating environment. In fact, by tuning the rejection criterion for the first classifier it is possible to reach the best trade-off between error and response time.

Specifically, we study a two-stage system where rejected patterns are forwarded to the nearest-neighbor classifier together with the $\text{top-}h$ ranking classes returned by the first classifier. Only patterns in the reference database belonging to the $\text{top-}h$ classes are used by the nearest-neighbor classifier to classify the rejected pattern. Moreover, to further speed-up the response time of the nearest-neighbor classifier, an editing of the nearest-neighbor reference database can be performed by collecting patterns rejected by the first-level classifier according to a narrower rejection criterion.

It is worth noting that this type of system is consistent with a view of the classification process which tries to conciliate a global approach with a local one. In fact, while the aim of global estimation is to estimate a function for all possible values of input, it is typically very hard, especially for multiclassification problems, to get a good estimate for inputs which are very close to the decision boundary. For these inputs, it is more effective to perform a local estimation (see [2], [33]) which focuses on a specific estimation point. Of course, the natural choice for performing local estimation is to use memory-based methods, such as nearest-neighbor.

The probability of error of the above system can be studied theoretically and it can be demonstrated that even if the two classifiers in the system are not optimal in a Bayes sense, the system as a whole, under specific conditions, may be optimal. Moreover, it is not difficult to formally relate the response time of the system to the rejection rate of the first classifier and to the other system parameters and, thus, to discuss the error-response time trade-off.

- N. Giusti is with Micronix Computer, Via dei Colombi, 2, 51016 Montecatini Terme (PT), Italy. E-mail: ngiusti@micronix.net.
- F. Masulli and A. Sperduti are with the Dipartimento di Informatica, Università di Pisa, Corso Italia 40, 56125 Pisa, Italy. E-mail: {masulli, perso}@di.unipi.it.

Manuscript received 4 Oct. 2000; revised 20 Sept. 2001; accepted 28 Nov. 2001.

Recommended for acceptance by M. Pietikainen.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 112942.

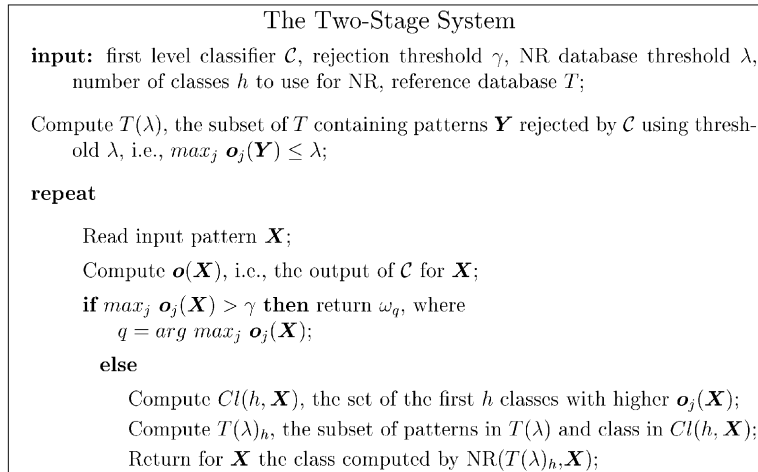


Fig. 1. The algorithm for the two-stage system.

We consider two specific instances of this scheme and we show that the obtained empirical results agree with the proposed framework, in the sense that the accuracy can be improved without significantly increase the average response time of the system.

The paper is organized as follows: The proposed two-stage classification system and its theoretical analysis (fully described in Appendix A) are presented in Section 2. Experimental results in this context are reported in Section 3, where two instances of the two-stage system are applied to handwritten digits recognition: in one instance, the first-level classifier is a fuzzy basis functions network (described in details in Appendix B), while in another instance a feed-forward neural network is used. The performances of the two systems for different instantiations of the parameters' values are computed and the experimental analysis shows that the two-stage approach can actually reach a higher performance than the ones obtained by its single component classifiers. Moreover, the analysis on the response time of the system can be used to select the different settings of parameters with the same response time, allowing the user to easily select the one which holds the best classification capability. Section 4 reports the conclusions.

2 A TWO-STAGE SYSTEM FOR MULTICATEGORY CLASSIFICATION TASKS

Let us to consider a multicategory classification framework, where an input pattern is represented as a random n -dimensional feature vector X belonging to one of N classes, $\omega_1, \dots, \omega_N$. In this context, we define a two-stage pattern recognition scheme consisting of a hierarchy made up by a multicategory (global) classifier \mathcal{C} with rejection, followed by a nearest-neighbor rule (NR) classifier working on the patterns rejected by \mathcal{C} . Specifically, a rejection rule is implemented as a *rejection threshold* on the level of the higher output, i.e., if no output of \mathcal{C} is greater than the threshold, the pattern is rejected from \mathcal{C} . The basic idea is that patterns rejected by the rejection rule with threshold value γ are then classified by the nearest-neighbor rule with reference database made up by patterns rejected by the same classifier \mathcal{C} , but using a threshold value $\lambda \geq \gamma$. By using the recognition threshold γ , \mathcal{C} classifies very fast most of the patterns with

small classification error, while a minority of patterns are forwarded to the NR for classification. Moreover, the quality of the NR database is controlled independently by the threshold λ . Of course, for rejected patterns, the recognition speed depends mainly on the dimension of the NR database. However, to speed-up the recognition time of the NR classifier, at classification time one can use an efficient *online editing strategy*, by considering for the NR only patterns belonging to classes that get the first h higher outputs by \mathcal{C} . The two-stage algorithm is defined in Fig. 1.

Let us denote the two-stage system with H . It is not difficult to show (see Appendix A) that the probability of error $e_{H,h}(\gamma)$ for the two-stage system, given a rejection threshold γ for the classifier at the first level, and selecting the top- h ranking classes for the second-level classifier, can be expressed as the sum of four positive terms

$$e_{H,h}(\gamma) = e_{bayes} + e_{\mathcal{C}}(\gamma) + e_{NR}(h) + e_{top-h}, \quad (1)$$

where e_{bayes} is the optimal Bayes error over the input domain, $e_{\mathcal{C}}(\gamma)$ is the error rate of the \mathcal{C} classifier for the given value of the rejection threshold γ , $e_{NR}(h)$ is the error induced by the application of the nearest-neighbor Rule on the patterns rejected by \mathcal{C} and depending on the value of h , and e_{top-h} is the error due the fact that the right class is not included (by \mathcal{C}) in the top- h classes. Notice that e_{top-h} is null if $h = N$. Moreover, under this condition, the two-stage system may reach the optimal Bayes error if and only if all the patterns which are not rejected by \mathcal{C} are classified correctly, while the rejected patterns are classified correctly by the nearest-neighbor rule. Of course, there is no guarantee that the two-stage system will reach the optimal Bayes error; however, the above decomposition of the error shows that, in principle, we do not loose the possibility to reach the optimal Bayes error even if both \mathcal{C} and NR are not optimal.

Finally, it should be noted that only under special conditions, i.e., using the k -nearest-neighbor rule with large values for k and an infinite amount of training examples, the optimal Bayes error can be reached: the NR by itself may only achieve an asymptotic error rate which is suboptimal (only about twice the Bayes error [6]). However, the use of a different second-level classifier turns out to be quite problematic. In fact, just training a global classifier on patterns

rejected by the first-level classifier does not lead to good performance since the complexity of the classification problem is exactly “coded” by the rejected patterns. This can be better understood when considering the support vectors of a (kernel) support vector machine [33]: The support vectors are the patterns closest to the decision boundary and training from scratch the classifier by just keeping these patterns will not change the outcome of learning.

Following this reasoning, it is clear that only a local classifier may be able to get a good performance. Given a testing pattern, an alternative to the use of NR would be to train a classifier with the training examples located in a small neighborhood around the testing pattern and then to apply the trained classifier to the testing pattern itself [2]. This approach, however, would excessively increase the response time of the system. Alternatively, a predefined decomposition of the input space in regions could be used to train once for all local classifiers assigned to each region [13]. This solution, however, could be inadequate since the a priori decomposition could turn out to be suboptimal.

2.1 Error-Response Time Trade-Off

For some applications, the response time of a system may be as important as the performance in generalization. The two-stage system is flexible enough to allow the balance of these two aspects by tuning the rejection threshold γ .

First of all, note that the classifier at level 1 has a constant response time, independently of the pattern in input. Let C_1 be the computational cost of running it on a pattern. On the other hand, the response time of the nearest-neighbor classifier using the top- h ranking selection rule depends on the number of prototypes stored for each of the top- h ranking classes. Let $\xi_i^h(\gamma) = P(\omega_i \in Cl(h)|\gamma)$, i.e., the probability of class ω_i being in the top- h ranking classes given a rejection threshold γ , then the expected number of prototypes $P_{2,h}(\gamma)$ used by the nearest-neighbor classifier is

$$P_{2,h}(\gamma) = h \sum_{i=1}^N c_i(\lambda) \xi_i^h(\gamma), \quad (2)$$

where $c_i(\lambda)$ is the number of prototypes of class ω_i in the reference database obtained by λ . Note that, if the prototypes are balanced across classes, i.e., $\forall i c_i(\lambda) = c_\lambda$, then $P_{2,h}(\gamma) = hc_\lambda$. The computational cost $C_{2,h}(\gamma)$ can then be defined as

$$C_{2,h}(\gamma) = \mathcal{I}(P_{2,h}(\gamma)), \quad (3)$$

where $\mathcal{I}(\cdot)$ is a function which depends on the way the nearest-neighbor classifier is implemented. Finally, the response time $C_{H,h}(\gamma)$ for the two-stage system is

$$C_H(\gamma) = C_1 + r(\gamma)C_{2,h}(\gamma), \quad (4)$$

where

$$r(\gamma) = \int_R p(X)dX. \quad (5)$$

Note that, in general, a local classifier needs much more time to respond than a global classifier. In our setting, the response time for the classifier at the second level can be reduced by using small values for both h and λ . This reduction in response time, however, is paid with a loss in generalization. Thus, the problem is to find a good trade-off

between $C_H(\gamma)$ and $e_{H,h}(\gamma)$. Since the relationship between the different components and parameters of the two-stage system are not linear, in this paper we have studied this problem from an experimental point of view.

3 EXPERIMENTAL RESULTS

As test-bed for our experiments we used the classification of handwritten digits. Specifically, we used a training set, a validation set, and a test set extracted from the NIST-3 database [8]. Both the training set and the validation set were made up of 10,000 associative pairs of segmented handwritten digits each, obtained from disjoint groups of writers, while the test set consisted of 5,000 independent digits.

The preprocessing of the digits included the following steps:

1. digit image extraction from the CD-ROM and normalization to a 32×32 binary matrix;
2. low-pass filtering in order to remove some small spots and holes from the image;
3. application of a shear transform to the digit image to straighten the axis joining the first upper-left point of the digit image to the last lower-right point;
4. image skeletonization by using a thinning algorithm;
5. finally, transformation of the digit representation into a 64-element vector, each vector element representing the number of black pixels contained in adjacent 4×4 squares (local counting).

It is worth noting that the resulting digit representation exhibits sufficient degrees of invariance to both scale and small image shifts or rotations.

3.1 Performance

The hardness of the classification problem was evaluated by studying the performance of the k -nearest-neighbor rule (k-NR) for different values of k (see Fig. 2). Due to sparseness of data, the best performance was obtained for $k=1$ (92.89 percent). In order to have a more extensive comparison against other classification techniques based on supervised training, we report in Table 1 the performance of the best fuzzy basis functions network (FBFN) (see [3] and Appendix B for details on the model) with 48, 12, and 10 hidden nodes we were able to obtain on the data. It must be pointed out that the FBFN holds universal approximation capabilities and under the same training conditions we adopted, it can approximate the Bayes optimal classifier [18]. Moreover, all the multilayer neural networks (NNs) we were able to train showed a slightly inferior performance with respect to FBFNs with a similar number of free parameters. As reference for NN, we will use a network with one hidden layer with 48 hidden units (NN₄₈), which achieved a performance of 93.42 percent on the test set.

In Table 1, we have also reported the results obtained on a simplified version of the FBFN (see Appendix B for details), namely, a ESFBFN with 20 hidden units. This network can be trained in less than half the time required by the FBFN₄₈ as well as its response time is less than half of the FBFN₄₈, while still preserving a similar performance in generalization.

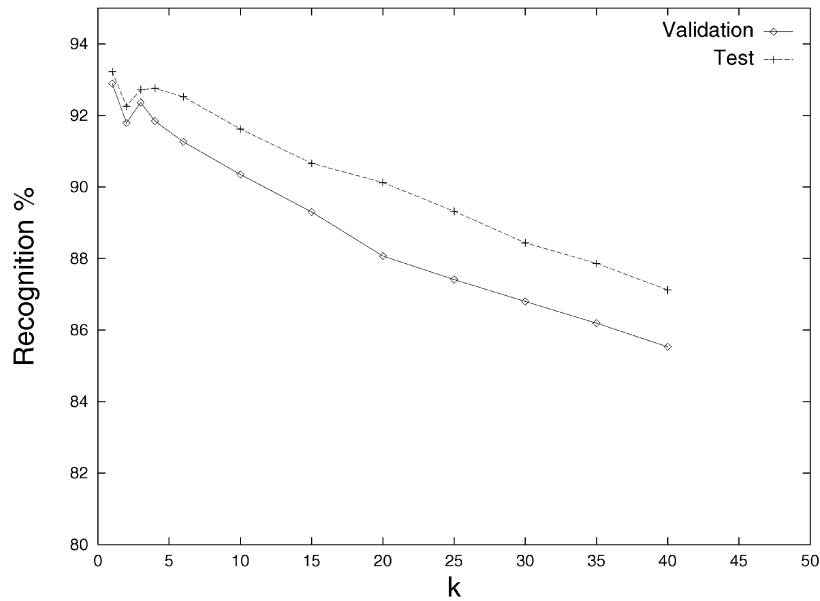


Fig. 2. Performance of the k -nearest-neighbor on the test and validation sets using as reference database the full training set.

TABLE 1

Comparison among FBF Networks (with 48, 12, and 10 Hidden Units) and a ESFBF Network with 20 Hidden Units (Two for Each Class), %S-Test Is the Success Rate on the Test Set, and the Epoch Duration is Measured on a Sun 10

MODEL	%S-Test	Epochs	Epochs Duration (sec)	Training Cost
FBFN ₄₈	94.09	13	1925	25025
FBFN ₁₂	92.26	36	307	11052
FBFN ₁₀	92.23	55	180	9900
ESFBFN ₂₀	93.80	250	49	12250

In order to avoid overfitting, the training for all the networks (including the neural network) was stopped using the validation set (early stopping).

Due to the better trade-off between response time and generalization performance, among the neuro-fuzzy models we selected the ESFBFN₂₀ as classifier at the first level of the first instance of the two-stage system. We also considered a system where the first classifier was the neural network NN₄₈ with normalized output. Moreover, accordingly with the results reported in Fig. 2, we decided to use $1 - NR$ as second-level classifier for both systems.

In Fig. 3, we have reported the distribution for classes of the reference database for the NR ($k = 1$) (to be used in the two-stage system) generated by ESFBFN₂₀ for different values of λ , while the effect of the rejection threshold on the training and test set for the ESFBFN₂₀ is shown in Fig. 4. Similar curves are obtained for NN₄₈.

Let us now turn to the performance of the whole two-stage systems. In Fig. 5, we have reported the test curves for $\gamma = 0.96$ and different values of λ for both systems. As expected, the performance of both systems improved with the dimension of the reference database for the NR (higher values of λ). Mixed results are instead obtained when considering the number of classes used for the NR (online editing): The system based on NN₄₈ got the worst performance with just two classes, while the system based on ESFBFN₂₀ got in this case its best performance. Thus, it seems that NN₄₈ is not able to get a very good estimate of the a posteriori probability of the first two best classes.

To assess the "best" values for γ and h (i.e., the number of classes to be used by the NR), we performed experiments with $\lambda = 1.0$ and different values for γ and h for both systems. The results obtained for the test set for both systems are reported in Fig. 6. On the sampled values for the parameters, the system based on NN₄₈ got the lowest error (5.08 percent) with $\gamma = 0.68$ and $h = 5$, while the system based on ESFBFN₂₀ got the lowest error (4.98 percent) with $\gamma = 0.88$ and $h = 2$.

In Fig. 7, we have reported the performance on the test set of both the NR and NN₄₈ on the set of patterns rejected by NN₄₈ for different values of γ . From these curves, it is clear that the NR outperforms NN₄₈ on patterns which are close to the decision boundary of NN₄₈. Similar results are obtained for ESFBFN₂₀.

3.2 Response Time

In order to study the response time of the system, we focus on the two-stage system where the first classifier is ESFBFN₂₀. For this system, we have estimated the quantities defined in (4), i.e., C_1 , $C_{2,h}$, and $r(\gamma)$. Both C_1 and $C_{2,h}$ have been estimated by considering the most relevant mathematical operations performed by \mathcal{C} and NR. We have considered additions, subtractions, multiplications, divisions, and the computation of the exponential function. The cost of performing 100 millions of each operation on a SUN 20 workstation has been computed experimentally. We obtained the same cost (say unitary cost) for addition, subtraction, and multiplication, while the division costed 2.208 and the computation of the

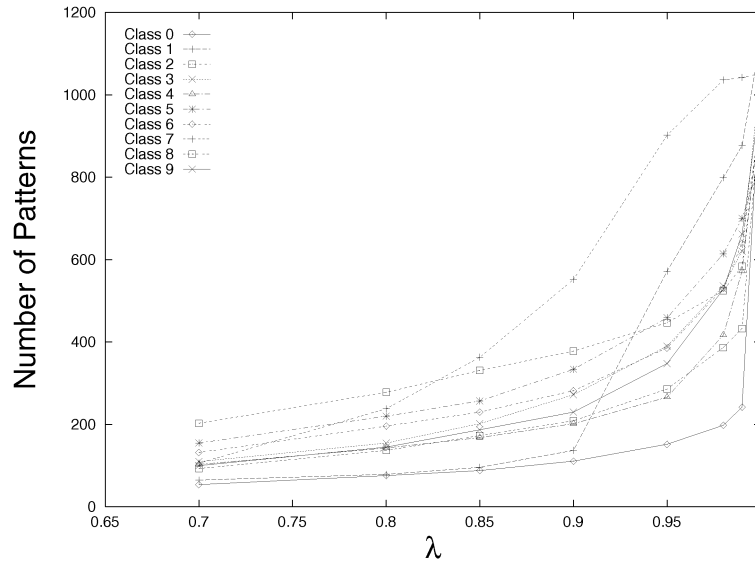


Fig. 3. Composition of the reference database for the NR ($k = 1$) for different classes and λ values by using ESFBFN₂₀. Note that the database for $\lambda = 1.0$ is equal to the training set.

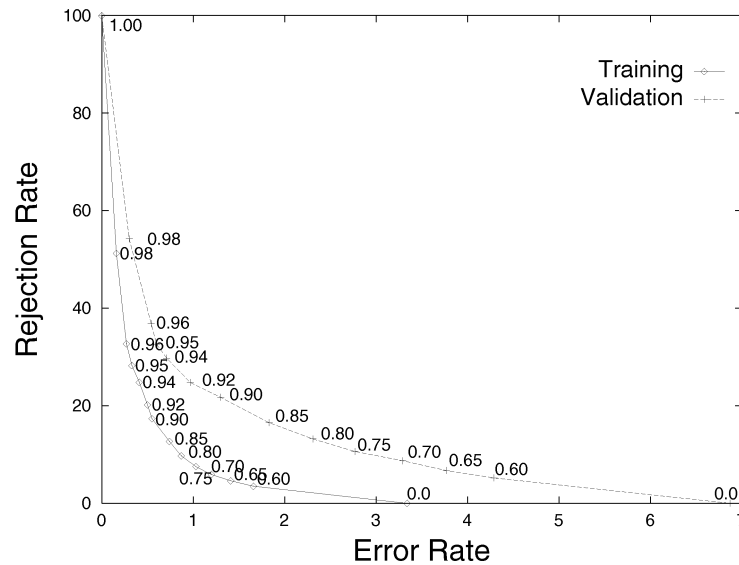


Fig. 4. Effect of the rejection threshold (reported near each experimental point) on the training and validation set for the ESFBFN₂₀. The numbers shown on the curves are the values of the γ threshold for which that error/rejection value was obtained.

exponential function 6.305. Using these weights, we estimated the computational cost for computing the output of \mathcal{C} and NR.¹ The rejection rate $r(\gamma)$ was experimentally computed over different values of γ close to the optimum value for the validation set, i.e., $\gamma = 0.87$, $\gamma = 0.88$, and $\gamma = 0.89$. The estimation of $C_H(\gamma)$ for different values of h and λ are reported in Fig. 8. As expected, as soon as h and λ increase, $C_H(\gamma)$ increases exponentially. This is particularly true with the increase of h . Moreover, the rate of increase is controlled by the value of γ . For each plot in Fig. 8, we have drawn on the cost surface curves with constant cost (i.e., $C_H(\gamma) = 100 \cdot 10^3$, $C_H(\gamma) = 200 \cdot 10^3$, and $C_H(\gamma) = 290 \cdot 10^3$). This was done to show that different couples of values for λ and h may end up to have the same computational cost and, thus, they have the same average response time.

1. No optimized algorithm for the NR has been considered.

If a bound on the average response time is given, these plots can be used to select the values of γ , λ , and h which are consistent with the target response time. Furthermore, the setting corresponding to the best performance can be chosen for the working system. In Fig. 9, we give an example of how this selection can be done. We have chosen a cost of 10^5 , which is used as cut point for the surfaces $C_H(\gamma)$ computed for sampled values of γ . The level curves obtained in this way are projected on the $\lambda - h$ plane and the performance of the two-stage system is evaluated on admissible points of the curves (i.e., $\lambda \geq \gamma$ and integer values for h) by using the validation set (Fig. 9a). The values for γ , λ , and h corresponding to the best performance (in this case, $\gamma = 0.8$, $\lambda = 0.988$, $h = 3$) are then used in the working system. From Fig. 9b, it can be noted that the selected values, as well as the performance, are not far from the optimal values for the test set, i.e., $\gamma = 0.85$, $\lambda = 0.983$, $h = 2$.

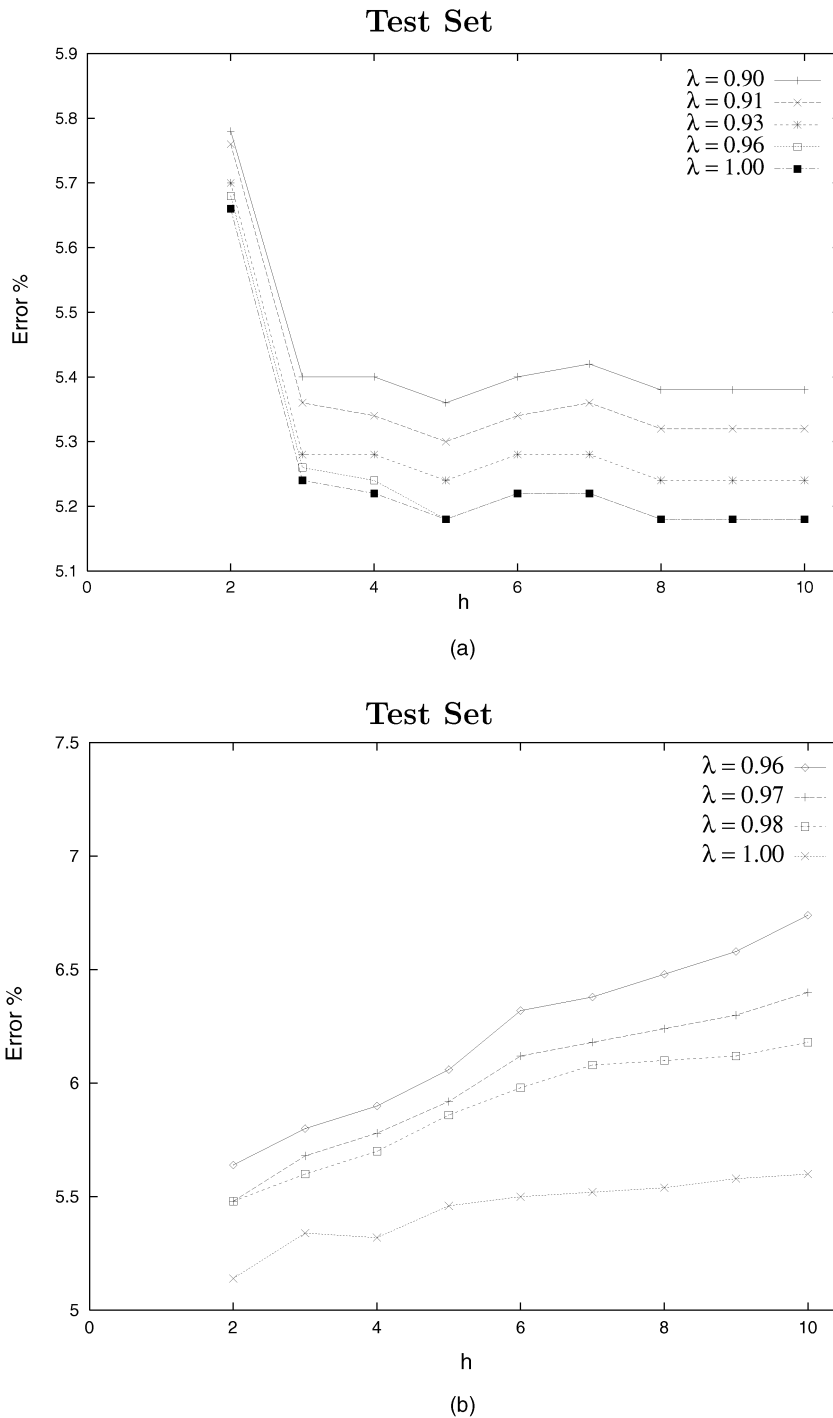


Fig. 5. Error curves for two different instances of the two-stage system for a fixed value of γ and different values of λ . (a) The first classifier is a NN with 48 hidden units ($\gamma = 0.65$). (b) The first classifier is a ESFBFN with 20 (hidden) units ($\gamma = 0.96$).

4 CONCLUSION

In the context of multicategory classification tasks, we have considered a two-stage system which combines a first-level global classifier with the nearest-neighbor Rule. This system is consistent with a view of the classification process which tries to conciliate a global approach with a local one in order to improve the trade-off between classification accuracy and response time. This trade-off is an important issue when considering classification tasks involving a high number of different classes.

For the proposed scheme, it is possible to theoretically relate the error rate of the system with the optimal Bayes error, showing that it is actually possible to reach the optimal Bayes error even if the two classifiers in the system are not optimal in a Bayes sense. Moreover, we formally related the expected average response time of the system to the rejection rate of the first classifier and to the other system parameters. This allows the tuning of the system parameters in order to reach the desired error-response time trade-off.

On two specific instances of the system and on a specific classification task, we have demonstrated that the

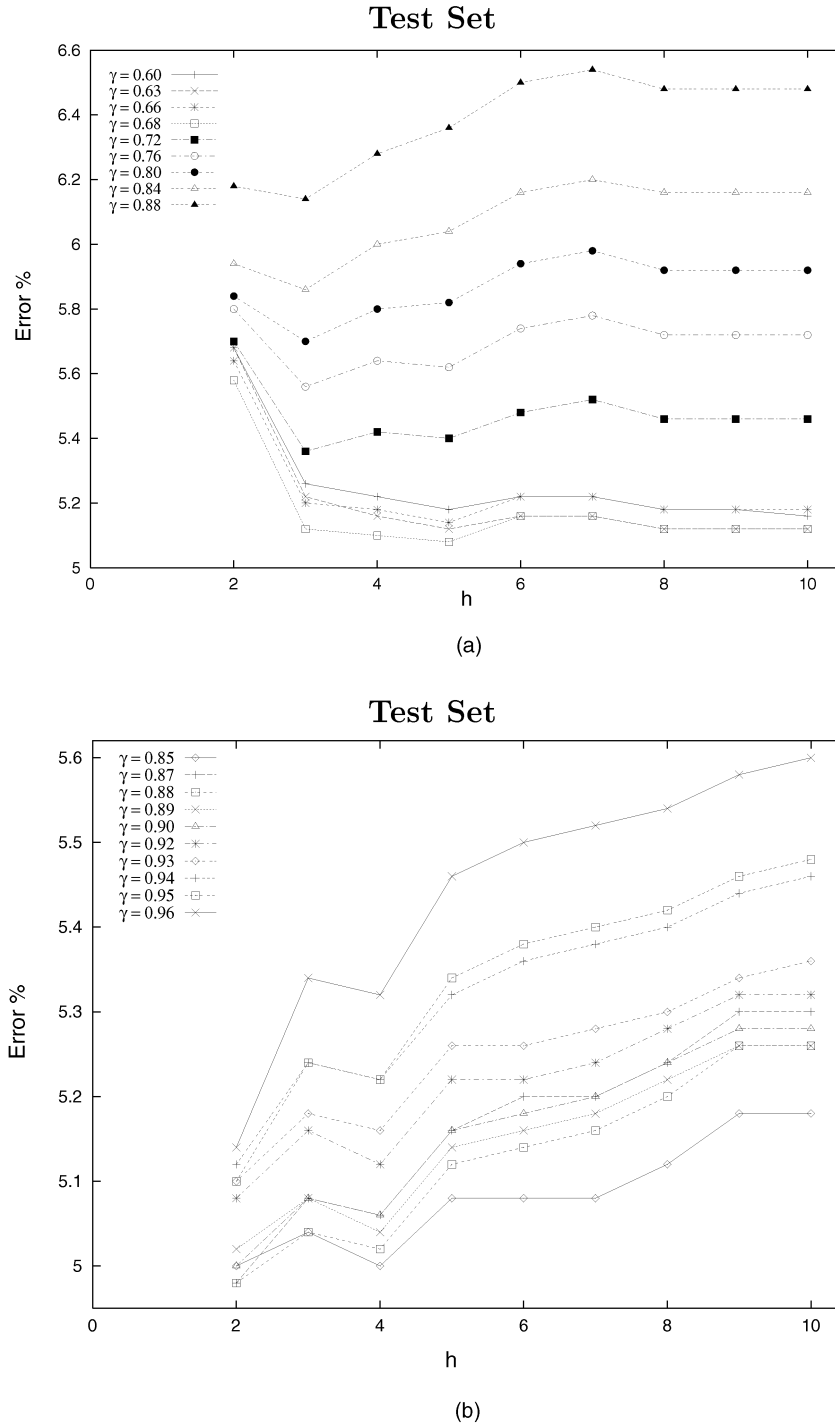


Fig. 6. Error curves for two different instances of the two-stage system for $\lambda = 1.0$ and different values of γ . (a) The first classifier is a NN with 48 hidden units. (b) The first classifier is a ESFBFN with 20 (hidden) units.

classification accuracy of the system can actually be higher than any of the single compounding classifiers. The system expected average response time is within² the range defined by the two compounding classifiers and it can be adjusted by tuning the system parameters.

We believe that for many other different instances of the system and many different classification tasks, a similar behavior, as predicted by the theoretical analysis, may be reproduced.

2. Actually, it is closer to the extreme defined by the first-level classifier.

APPENDIX A

THEORETICAL ANALYSIS OF THE TWO-STAGE SYSTEM

In this appendix, we perform a theoretical analysis of the two-stage system error, showing in detail the decomposition of the error reported in (1).

Let π_{ω_i} denote the a priori probability of observing class ω_i , while the *posterior* probability of class ω_i given an input vector X is denoted as

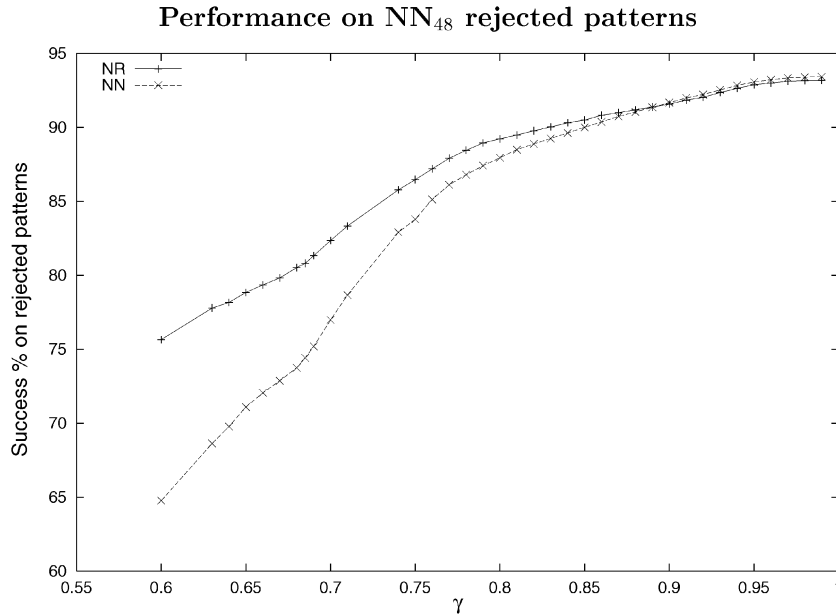


Fig. 7. Performance of the NR (with full training set) and NN_{48} on the patterns rejected by the NN_{48} for different values of γ .

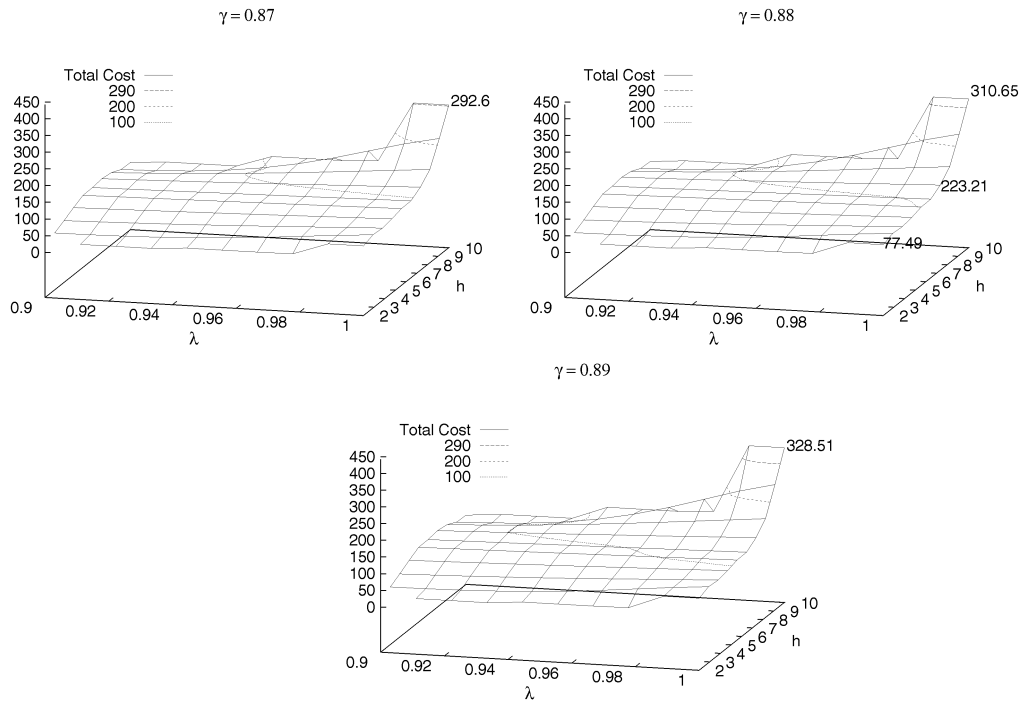


Fig. 8. Plots representing an estimate of $C_H(\gamma)$ for different values of γ . $C_H(\gamma)$ is plotted for different values of h and λ . Costs must be multiplied for 10^3 .

$$P(\omega_i | \mathbf{X}) = \frac{p(\mathbf{X} | \omega_i)\pi_{\omega_i}}{p(\mathbf{X})}, \quad (6)$$

where $p(\mathbf{X} | \omega_i)$ is the ω_i class conditional probability density function and $p(\mathbf{X})$ is the input vector probability density function.

For our aims, it is important to sort the posterior probabilities $P(\omega_i | \mathbf{X})$ in decreasing order. With

$$P_1(\mathbf{X}) \geq P_2(\mathbf{X}) \geq \dots \geq P_N(\mathbf{X}), \quad (7)$$

we denote the ordered sequence of posterior probabilities, where

$$P_1(\mathbf{X}) = \max_{j \in [1, \dots, N]} P(\omega_j | \mathbf{X}), \quad (8)$$

and so on. In the two-stage system, the output $o(\mathbf{X})$ of the classifier \mathcal{C} is actually providing an approximation of the posterior probabilities. However, since $o(\mathbf{X})$ is only an approximation of the true posterior probabilities, the order induced by $o(\mathbf{X})$ over the classes will, in general, be different from the one induced by the true posterior probabilities. Consequently, let

$$F_1(\mathbf{X}), F_2(\mathbf{X}), \dots, F_N(\mathbf{X}) \quad (9)$$

be sequence (7) reordered according to $o(\mathbf{X})$, i.e.,

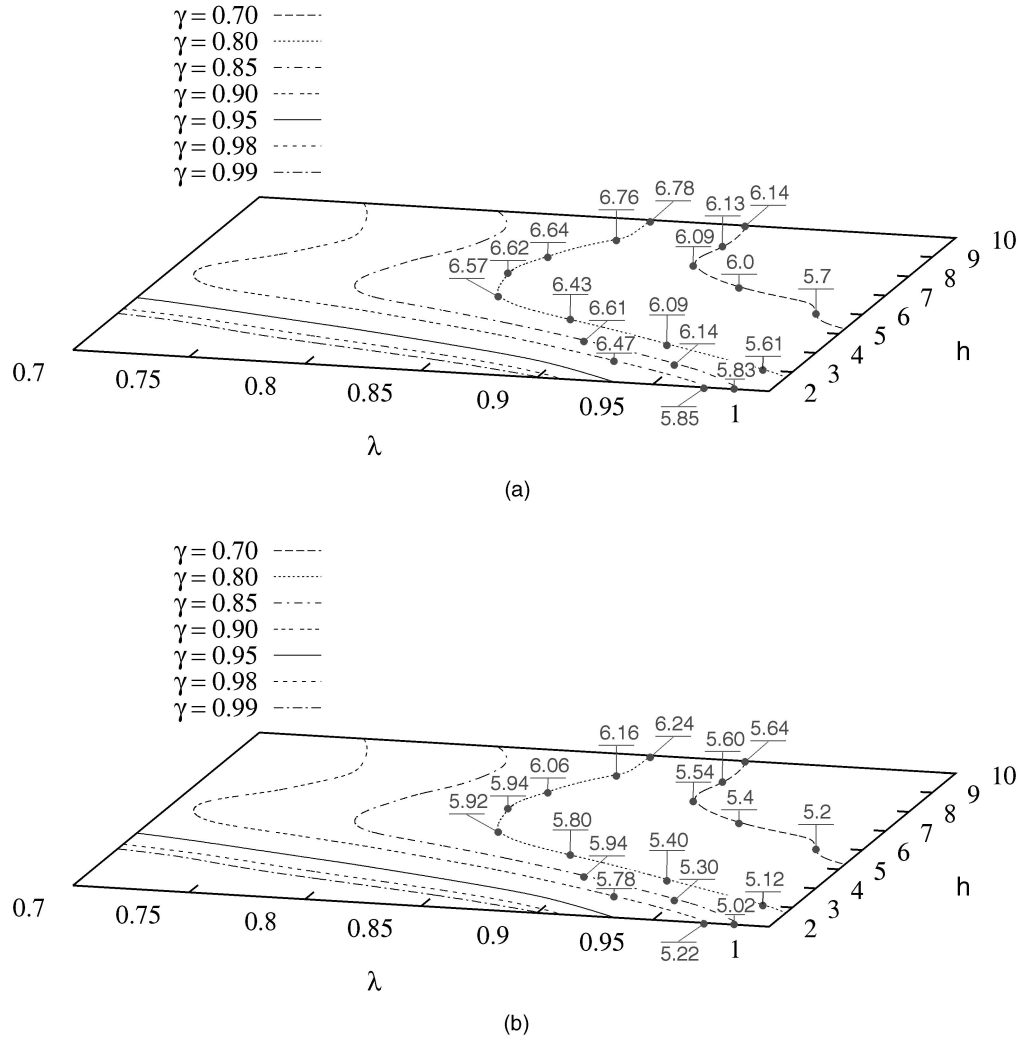


Fig. 9. Level curves of the estimated $C_H(\gamma)$ for different values of γ given a fixed cost of 10^5 . The performance obtained by a two-stage system with a ESFBFN with 20 (hidden) units on (a) the validation set and (b) the test set, for admissible points on the curves (i.e., $\lambda \geq \gamma$ and integer values for h), are shown. Given a cost (e.g., 10^5), the corresponding plot on the validation set (e.g., (a)) can be used to decide the best setting for λ and h (for this example, $\gamma = 0.8$, $\lambda = 0.988$, $h = 3$, which is almost optimal with respect to the test set: see plot (b)). Notice that plot (b) is consistent with plot (a).

$$F_1(\mathbf{X}) = P(\omega_k | \mathbf{X}) \quad (10)$$

with

$$k = \operatorname{argmax}_{j \in [1, \dots, N]} o_j(\mathbf{X}), \quad (11)$$

and so on. We can now define the discrepancy values induced by the classifier \mathcal{C} as

$$\Delta^F(\mathbf{X}) = P_1(\mathbf{X}) - F_1(\mathbf{X}). \quad (12)$$

The same definitions can be similarly devised for the classifier at the second level of the two-stage system (nearest-neighbor):

$$G_1(\mathbf{X}), G_2(\mathbf{X}), \dots, G_N(\mathbf{X}), \quad (13)$$

and

$$\Delta^G(\mathbf{X}) = P_1(\mathbf{X}) - G_1(\mathbf{X}). \quad (14)$$

Note that $\Delta^F(\mathbf{X}) \geq 0$ and $\Delta^G(\mathbf{X}) \geq 0$.

Let us denote the two-stage system with H . The probability of error $e_{H,h}(\gamma)$ for the two-stage system, when using a rejection threshold γ for the classifier at the first

level, and selecting the top- h ranking classes for the second-level classifier, can be expressed as

$$e_{H,h}(\gamma) = e_1(\gamma) + e_{2,top-h}(\gamma), \quad (15)$$

where $e_1(\gamma)$ is the error rate of the \mathcal{C} classifier, and $e_{2,top-h}(\gamma)$ is the error induced by the application of the second-level classifier selecting only the top- h ranking classes suggested by the \mathcal{C} classifier.

The first-level error rate can be written as

$$e_1(\gamma) = \int_A (1 - F_1(\mathbf{X})) p(\mathbf{X}) d\mathbf{X}, \quad (16)$$

where A is the input subspace accepted for classification by the classifier, i.e., $A = \{\mathbf{X} \mid \max_{j \in [1, \dots, N]} o_j(\mathbf{X}) \geq \gamma\}$. By using the discrepancy values, $e_1(\gamma)$ can be cast in the following form

$$e_1(\gamma) = \int_A (1 - F_1(\mathbf{X}) + P_1(\mathbf{X}) - P_1(\mathbf{X})) p(\mathbf{X}) d\mathbf{X} \quad (17)$$

$$= \int_A (1 - P_1(\mathbf{X})) p(\mathbf{X}) d\mathbf{X} + \int_{A \cap I^F} \Delta^F(\mathbf{X}) p(\mathbf{X}) d\mathbf{X}, \quad (18)$$

where $I^F = \{X \mid \Delta^F(X) \neq 0\}$. On the other side, $e_{2,top-h}(\gamma)$ can be written as

$$e_{2,top-h}(\gamma) = \int_R (1 - G_1(X) \sum_{j=1}^h F_j(X)) p(X) dX, \quad (19)$$

where R is the input subspace rejected by the first-level classifier, i.e., $R = \{X \mid \max_{j \in [1, \dots, N]} o_j(X) < \gamma\}$. After some algebra it turns out that

$$\begin{aligned} e_{2,top-h}(\gamma) &= \int_R (1 - P_1(X)) p(X) dX + \int_{R \cap I^G} \Delta^G(X) p(X) dX \\ &+ \int_R G_1(X) \sum_{j=h+1}^N F_j(X) p(X) dX, \end{aligned} \quad (20)$$

where $I^G = \{X \mid \Delta^G(X) \neq 0\}$.

Thus, the error rate for the two-stage system can be written as

$$\begin{aligned} e_{H,h}(\gamma) &= \\ e_{bayes} &+ \int_{A \cap I^F} \Delta^F(X) p(X) dX + \int_{R \cap I^G} \Delta^G(X) p(X) dX \\ &+ \int_R G_1(X) \sum_{j=h+1}^N F_j(X) p(X) dX, \end{aligned} \quad (21)$$

where e_{bayes} is the optimal Bayes error over the input domain, and all the remaining terms are nonnegative. Note that, in order for the two-stage system to reach the optimal Bayes error, we must have $h = N$, $A \cap I^F = \emptyset$, and $R \cap I^G = \emptyset$. So, even if the classifiers at the first and second level are not optimal in the Bayes sense,³ the two-stage system can still approach the optimal error rate. What is important is that the classifier at the first level must misclassify only vectors in R , while the classifier at the second level must misclassify only vectors in A .

Of course, it is not reasonable to expect these kind of behavior from the classifiers in the two-stage system. However, it must be noted that the error of the classifier at the first level decreases with the increase of γ . In fact, when considering a bayesian classifier, the adopted rejection rule⁴ guarantees that it is possible to express the error rate directly as a function of the reject rate via the following Stieltjes integral [4]

$$e_1(\bar{t}) = - \int_0^{\bar{t}} t d\rho(t), \quad (22)$$

where $\gamma = 1 - t$,

$$\rho(t) = \int_{R_{bayes}} p(X) dX, \quad (23)$$

and $R_{bayes} = \{X \mid \max_{j \in [1, \dots, N]} P_j(X) < 1 - t\}$.

The classifier at the first level will have, in general, a small error probability, which means that the set $A \cap I^F$ is small in size. Moreover, the pattern rejected by it will be

located, in general, close to the boundaries of the classes. As a consequence, the rejected patterns are expected to be rather sparse, and a further classification by using a global classifier is not going to return a satisfactory classification. In order to correctly classify these patterns, it is more productive to use a local classifier [33], such as the nearest-neighbor classifier. Because of the above considerations, the classifier \mathcal{C} will tend to minimize the term

$$e_C(\gamma) = \int_{A \cap I^F} \Delta^F(X) p(X) dX, \quad (24)$$

while the classifier at the second level, because of its locality, should reduce considerably the term

$$e_{NR}(h) = \int_{R \cap I^G} \Delta^G(X) p(X) dX. \quad (25)$$

The remaining term

$$e_{top-h} = \int_R G_1(X) \sum_{j=h+1}^N F_j(X) p(X) dX, \quad (26)$$

is minimized by setting $h = N$; however, for speeding up the response time of the system, it may be convenient to have $h < N$.

APPENDIX B

THE EXTENDED SIMPLIFIED FUZZY BASIS FUNCTION NETWORK

A Fuzzy Logic System with *singleton* fuzzification, *max-product* composition, *product inference*, and *height* defuzzification can be represented as [19]

$$y = f(\mathbf{x}) = \sum_{l=1}^M \bar{y}^l \phi_l(\mathbf{x}) \quad (27)$$

with

$$\phi_l(\mathbf{x}) = \frac{\prod_{i=1}^p \mu_{F_i^l}(x_i)}{\sum_{l=1}^M \prod_{i=1}^p \mu_{F_i^l}(x_i)}, \quad (28)$$

where \bar{y}^l denotes the center of gravity of the output fuzzy set, and $\phi_l(\mathbf{x})$, $l = 1, 2, \dots, M$, are called *fuzzy basis functions*. We can refer to those FLS as *fuzzy basis expansions* or *networks of fuzzy basis functions* (FBF network).⁵

The relationships between fuzzy basis expansions and other basis functions have been extensively studied in [12]. It is worth noting that the FLS with universal function property studied by Mendel and Wang [35], [34] (i.e., a singleton FLS using product inference, product implication, Gaussian membership, and height defuzzification) can be rewritten as a FBF network expansion.

Here, we are interested in a neuro-fuzzy logic system based on a multi-input-multi-output (MIMO) version of this FBF network. Specifically, if there are K units in the input layer, J fuzzy inference rules and I outputs, the rule activations can be expressed as:

5. In [19], fuzzy basis expansions for FLS with nonsingleton fuzzification are also introduced.

3. Note that I^F or I^G are not required to be empty.

4. The rejection rule we adopted corresponds to the Chow rule [4]. Chow has shown that this rule is optimal in the sense that for a given error rate (error probability) the rejection rate (reject probability) is minimized.

$$r_j = \prod_k \mu_{jk}(x_k) \quad (29)$$

$$\mu_{jk}(x_k) = \exp\left(-\frac{(x_k - m_{jk})^2}{2\sigma_{jk}^2}\right) \quad (30)$$

$$y_i = \frac{\sum_j r_j \bar{y}_{ij}}{\sum_j r_j} = \sum_j \bar{y}_{ij} \phi_j(\mathbf{x}) \quad (31)$$

$$\phi_j = \frac{\prod_k \mu_{jk}(x_k)}{\sum_j \prod_k \mu_{jk}(x_k)}, \quad (32)$$

where the quantity $\mu_{jk}(x_k)$ represents the value of the membership function of the component x_k of the input vector for the j th rule, m_{jk} and σ_{jk}^2 are the means and the variances of the Gaussian membership functions, y_i are the values of the output units, \bar{y}_{ij} is the center of gravity of the output fuzzy membership function of the j th rule associated with the output y_i , and ϕ_j is the fuzzy basis function associated to rule j , representing its normalized activation.⁶

The FBF network can be regarded as a fuzzy system mapped on a network of RBF. The FBF network can be identified both by exploiting the linguistic knowledge available (*structure identification problem*) [15] and by using the information contained in a data set (*parameter estimation problem*) [15]. Learning rules based on Gradient Descent technique are discussed, e.g., in [35].

For pattern recognition applications, from this FBF network a *Simplified FBF network* (SFBF network) can be obtained by assuming, in accordance with *rule specialization* [1]:

$$\bar{y}_{ij} \equiv \delta_{ij} = \begin{cases} 1 & \text{if rule } j \text{ is associated to class } i, \\ 0 & \text{otherwise.} \end{cases} \quad (33)$$

This assumption leads to both a system with as many units as classes and a strong simplification of the learning formulas that become:

$$\Delta m_{jk} = \eta_m \phi_j U_{ij} [x_k - m_{jk}] / \sigma_{jk}^2 \quad (34)$$

$$\Delta \sigma_{jk} = \eta_\sigma \phi_j U_{ij} [x_k - m_{jk}]^2 / \sigma_{jk}^3 \quad (35)$$

with

$$U_{ij} = \begin{cases} (y_i - 1)^2 & \text{if } j = i \\ y_i^2 - y_i & \text{if } j \neq i. \end{cases} \quad (36)$$

It is worth noting that, from (33) and the form of the defuzzifier, $y \in (0, 1)$ follows, and, consequently,

$$U_{ij} = \begin{cases} \geq 0 & \text{if } j = i \\ \leq 0 & \text{if } j \neq i. \end{cases} \quad (37)$$

holds.

Therefore, the learning rules of the SFBF network are competitive. During training, the means of the Gaussian membership functions of each rule move toward the patterns of the class associated to that rule and escape from patterns belonging to other classes. At the same time, sigmas of Gaussian membership functions of each rule grow in order to increase the value of the membership function for patterns of the class associated to that rule or shrink in order to reduce the value of the Gaussian membership function for patterns

belonging to other classes. From a probabilistic point of view, the SFBF can be seen as a mixture of Gaussians with diagonal covariance matrices.

Of course, since this system must have as many units as classification classes, it cannot be used for complex classification tasks. To remove this constraint, a new level of competition among units can be introduced. The new defined network possesses n_j units associated to each class j , for a total of $J = \sum_{j=1}^I n_j$ units. During learning, the output of each unit is computed and the best unit for each class is selected, i.e., for each class j , the unit $i_j^* \in Idx_j = \{1, \dots, n_j\}$ such that $i_j^* = \arg \max_{i \in Idx_j} \{\phi_i\}$ is selected. In that way, the number of selected units is equal to the number of classes and the learning rules of the SFBF network can be applied. Thus, at each learning step, only the selected rules have the weights changed. During the operational phase, the input pattern is classified by the class label associated with the unit having maximum activity.

ACKNOWLEDGMENTS

This work was supported by grants from CNR, INFM, and MURST. The authors would like to thank Stefano Rovetta for helpful suggestions.

REFERENCES

- [1] D. Alfonso, F. Masulli, and A. Sperduti, "Competitive Learning in a Classifier Based on an Adaptive Fuzzy System," *Proc. Int'l ICSC Symp. Industrial Intelligent Automation (IIA '96) and Soft Computing (SOCO '96)*, P.G. Anderson and K. Warwick, eds., pp. 2-8, 1996.
- [2] L. Bottou and V. Vapnik, "Local Learning Algorithms," *Neural Computation*, vol. 4, no. 6, pp. 888-901, 1992.
- [3] F. Casalino, F. Masulli, and A. Sperduti, "Rule Specialization in Networks of Fuzzy Basis Functions," *Intelligent Automation and Soft Computing*, vol. 4, pp. 73-82, 1998.
- [4] C.K. Chow, "On Optimum Recognition Error and Reject Trade-Off," *IEEE Trans. Information Theory*, vol. 16, pp. 41-46, 1970.
- [5] W.B. Croft and L.S. Larkey, "Combining Classifiers in Text Categorization," *Proc. 19th Ann. Int'l Conf. Research and Development in Information Retrieval (SIGIR 96)*, pp. 289-297, 1996.
- [6] R.O. Duda and P.E. Hart, *Pattern Classification and Scene Analysis*. New York: Wiley, 1973.
- [7] C. Furlanello, D. Giuliani, E. Trentin, and S. Merler, "Speaker Normalization and Model Selection of Combined Neural Networks," *Connection Science*, vol. 9, no. 1, pp. 31-50, 1997.
- [8] M.D. Garris and R.A. Wilkinson, *NIST Special Database3 Handwritten Segmented Characters*. Gaithersburg, Md.: Nat'l Inst. of Standard and Technology, 1992.
- [9] S. Gutta and H. Wechsler, "Gender Classification of Human Faces Using Hybrid Classifier Systems," *Proc. Int'l Conf. Neural Networks*, vol. 3, pp. 1353-1358, 1997.
- [10] S. Hashem, "Effects of Collinearity on Combining Neural Networks," *Connection Science*, vol. 8, nos. 3 and 4, pp. 315-336, 1996.
- [11] D. Jimenez and N. Walsh, "Dynamically Weighted Ensemble Neural Networks for Classification," *Proc. 1998 Int'l Joint Conf. Neural Networks*, pp. 753-758, 1998.
- [12] H.M. Kim and J.M. Mendel, "Fuzzy Basis Functions: Comparisons with Other Basis Functions," *IEEE Trans. Fuzzy Systems*, vol. 3, pp. 158-168, 1995.
- [13] S. Knerr and A. Sperduti, "Rejection Driven Hierarchy of Neural Network Classifiers," *Proc. Int'l Symp. Nonlinear Theory and Its Applications '93*, pp. 957-961, 1993.
- [14] M.W. Kurzynski, "On the Multistage Bayes Classifier," *Pattern Recognition*, vol. 22, no. 4, pp. 355-365, 1988.
- [15] C.C. Lee, "Fuzzy Logic in Control Systems: Fuzzy Logic Controller, I," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 20, pp. 404-418, 1990.

6. Without loss of generality, we could assume that the fuzzy membership functions are singletons ($\bar{y}_{ij} \equiv s_{ij}$).

- [16] Y. Liu and X. Yao, "A Cooperative Ensemble Learning System," *Proc. 1998 IEEE Int'l Joint Conf. Neural Networks (IJCNN '98)*, pp. 2202-2207, 1998.
- [17] R. Maclin and D. Opitz, "An Empirical Evaluation of Bagging and Boosting," *Proc. 14th Nat'l Conf. Artificial Intelligence and Proc. Ninth Innovative Applications of Artificial Intelligence Conf. (AAAI '97/IAAI '97)*, pp. 546-551, pp. 27-31, July 1997.
- [18] F. Masulli, F. Casalino, and F. Vannucci, "Bayesian Properties and Performances of Adaptive Fuzzy Systems in Pattern Recognition Problems," *Proc. European Conf. Artificial Neural Networks, (ICANN '94)* M. Marinaro and P.G. Morasso, eds., pp. 189-192, 1994.
- [19] J.M. Mendel, "Fuzzy Logic Systems for Engineering: A Tutorial," *Proc. IEEE*, vol. 83, pp. 345-377, 1995.
- [20] D.W. Opitz and J.W. Shavlik, "Actively Searching for an Effective Neural Network Ensemble," *Connection Science*, vol. 8, nos. 3 and 4, 1996.
- [21] B. Parmanto, P.W. Munro, and H.R. Doyle, "Reducing Variance of Committee Prediction with Resampling Techniques," *Connection Science*, vol. 8, nos. 3 and 4, 1996.
- [22] P. Poddar and P.V.S. Rao, "Hierarchical Ensemble of Neural Networks," *Proc. Int'l Conf. Neural Networks*, vol. 1, pp. 287-292, 1993.
- [23] P. Pudil, J. Novovicova, S. Blaha, and J. Kittler, "Multistage Pattern Recognition with Reject Option," *Proc. 11th IAPR Int'l Conf. Pattern Recognition*, vol. 2, pp. 92-95, 1992.
- [24] C. Rodriguez, J. Muguerza, M. Navarro, A. Zarate, J.I. Martin, and J.M. Perez, "A Two-Stage Classifier for Broken and Blurred Digits in Forms," *Proc. 14th Int'l Conf. Pattern Recognition*, vol. 2, pp. 1101-1105, 1998.
- [25] J. Rokui and H. Shimodaira, "Multistage Building Learning Based on Misclassification Measure," *Proc. Int'l Conf. Artificial Neural Networks*, pp. 221-226, 1999.
- [26] B.E. Rosen, "Ensemble Learning Using Decorrelated Neural Networks," *Connection Science*, vol. 8, nos. 3 and 4, pp. 373-384, 1996.
- [27] R.F. Schapire and Y. Singer, "Improved Boosting Algorithms Using Confidence-Rated Predictions," *Proc. 11th Ann. Conf. Computational Learning Theory (COLT '98)*, pp. 80-91, July 1998.
- [28] A.J.C. Sharkey, "On Combining Artificial Neural Nets," *Connection Science*, vol. 8, nos. 3 and 4, pp. 299-314, 1996.
- [29] A.J.C. Sharkey, "Modularity, Combining and Artificial Neural Nets," *Connection Science*, vol. 9, no. 1, pp. 3-10, 1997.
- [30] S. Simon, H.A. Kestler, A. Baune, F. Schwenker, and G. Palm, "Object Classification with Simple Visual Attention and a Hierarchical Neural Network for Subsymbolic-Symbolic Coupling," *Proc. IEEE Int'l Symp. Computational Intelligence in Robotics and Automation*, pp. 244-249, 1999.
- [31] S.K. Tso, X.P. Gu, Q.Y. Zeng, and K.L. Lo, "Input Space Decomposition and Multilevel Classification Approach for ANN-Based Transient Security Assessment," *Proc. Fourth Int'l Conf. Advances in Power System Control, Operation and Management*, vol. 2, pp. 499-504, 1997.
- [32] K. Tumer and J. Ghosh, "Error Correlation and Error Reduction in Ensemble Classifiers," *Connection Science*, vol. 8, nos. 3 and 4, 1996.
- [33] V.N. Vapnik, *Statistical Learning Theory*. Wiley & Sons, 1998.
- [34] L. Wang and J.M. Mendel, "Fuzzy Basis Functions, Universal Approximation, and Orthogonal Least-Squares Learning," *IEEE Trans. Neural Networks*, vol. 5, pp. 807-814, 1992.
- [35] L. Wang and J.M. Mendel, "Generating Fuzzy Rules by Learning from Examples," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 22, pp. 1414-1427, 1992.



Francesco Masulli received the Laurea degree in physics from the University of Genova in 1976. He is an associate professor of computer science with the University of Pisa, Italy. After serving in the military, he was a researcher with the Italian National Institute for Nuclear Physics (1978-1979), the Ansaldo Automazione Co. (1979-1983), and he was an assistant professor with the University of Genova (1983-2001). He was also on leave as a visiting scientist at the University of Nijmegen, Holland, in 1983, and at the International Computer Science Institute in Berkeley, California in 1991, 1993, and 1994. He has authored or coauthored more than 100 scientific papers on machine learning, neural networks, fuzzy systems, and ensemble methods. He has coedited three books and two special issues of scientific journals on those subjects. He serves as an associate editor for the international journal *Intelligent Automation and Soft Computing* and also serves as a cochair for the 2002 Course of the International School on Neural Networks E.R. Caianiello on "Ensemble Methods in Learning Machines." His previous duties include chairing the Conference of the International Graphonomics Society (IGS) in 1997 and the Symposium on Soft Computing SOCO, in 1999. He is member of the IEEE-Neural Network Council (Italian R.I.G.), an affiliate member of the Berkeley Initiative on Soft Computing (BISC), and a board member of the Italian Neural Network Society (SIREN) and of the SIG Italy of the International Neural Network Society (INNS). He is a member of the IEEE and the IEEE Computer Society.



Alessandro Sperduti received the Laurea and doctoral degrees in computer science from the University of Pisa, Italy, in 1988 and 1993, respectively. In 1993, he spent a period of time at the International Computer Science Institute, Berkeley, supported by a postdoctoral fellowship. In 1994, he returned to the Computer Science Department at the University of Pisa, where he was assistant professor and where he presently is an associate professor. His research interests include pattern recognition, image processing, neural networks, and hybrid systems. In the field of hybrid systems, his work has focused on the integration of symbolic and connectionist systems. He contributed to the organization of several workshops on this subject and he served also in the program committee of conferences on neural networks. He has authored more than 80 refereed papers mainly in the areas of neural networks, fuzzy systems, pattern recognition, and image processing. Moreover, he has given several tutorials within international schools and conferences, such as the International Joint Conference on Artificial Intelligence 1997 and 1999. He has acted as guest co-editor of the *IEEE Transactions on Knowledge and Data Engineering* for a special issue on connectionist models for learning in structured domains. He is a member of the executive board of the Italian Neural Network Society (SIREN) and the European Neural Networks Council. He is a member of the IEEE Computer Society.

► For more information on this or any other computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.



Nicola Giusti received the Laurea degree in computer science from the University of Pisa in 1996. He has been involved in the definition and developing of competitive and hybrid neuro-fuzzy models for supervised classification. He currently is with Micronix Computer S.p.A. as application consultant in networking (LAN/WAN).