

F. Masulli · G. Valentini

Effectiveness of error correcting output coding methods in ensemble and monolithic learning machines

Received: 3 November 2000 / Accepted: 10 February 2003

© Springer-Verlag London Limited 2003

Abstract Error Correcting Output Coding (ECOC) methods for multiclass classification present several open problems ranging from the trade-off between their error recovering capabilities and the learnability of the induced dichotomies to the selection of proper base learners and to the design of well-separated codes for a given multiclass problem. We experimentally analyse some of the main factors affecting the effectiveness of ECOC methods. We show that the architecture of ECOC learning machines influences the accuracy of the ECOC classifier, highlighting that ensembles of parallel and independent dichotomic Multi-Layer Perceptrons are well-suited to implement ECOC methods. We quantitatively evaluate the dependence among codeword bit errors using mutual information based measures, experimentally showing that a low dependence enhances the generalisation capabilities of ECOC. Moreover we show that the proper selection of the base learner and the decoding function of the reconstruction stage significantly affects the performance of the ECOC ensemble. The analysis of the relationships between the error recovering power, the accuracy of the base learners, and the dependence among codeword bits show that all these factors concur to the effectiveness of ECOC methods in a not straightforward way, very likely dependent on the distribution and complexity of the data.

Keywords Coding · Classification problems · ECOC ensemble · Ensemble of learning machines · Error correcting output

Francesco Masulli (✉)
INFN – Istituto Nazionale per la Fisica della Materia, via
Dodecaneso 33, 16146 Genova, Italy and
Dipartimento di Informatica, Università di Pisa, Via Buonozzoti,
2 56125, Pisa, Italy.
E-mail: masulli@di.unipi.it

Giorgio Valentini
INFN – Istituto Nazionale per la Fisica della Materia, via
Dodecaneso 33, 16146 Genova, Italy and
DSI, Dipartimento di Scienze dell' Informazione, Università
degli studi di Milano, via Comelico 39, Milano, Italy.
E-mail: valentini@dsi.unimi.it

Introduction

In the last decade several methods for constructing ensembles of learning machines have been developed [1]. These methods encompass a wide range of technique such as ensemble averaging [2,3], where the outputs of different predictors are linearly combined to produce an overall output; boosting [4,5,6] and bagging [7], where the same learning algorithms are applied to different subsets of the training sets; mixture of experts [8,9], where the outputs of the different predictors are non linearly combined through a gating network; ensemble constructed by subsets of input features [10], where each predictor selects a group of the input features.

This paper focuses on Error Correcting Output Coding (ECOC) decomposition methods [11–15], and in particular on the factors affecting the effectiveness of these ensemble methods.

Error correcting output codes [16], originally proposed to enhance the reliability of the transmission of binary signals through a noisy channel [17], have been successfully used in the framework of decomposition methods for multiclass classification problems to improve the generalisation capabilities of learning machines. By this approach an overall classification problem is decomposed into a set of simpler dichotomic subtasks, through a manipulation of the output targets assigned to each class. Dietterich and Bakiri [11,12] demonstrated that ECOC can achieve better performances than classification methods based on distributed output codes [18]. In fact, using *codewords* for coding the classes suggest the introduction of codes with error recovering abilities. Kong and Dietterich showed that ECOC techniques can also provide class probability informations, through the solution of an over-constrained system of linear equation [14]. An interesting extension of this approach presented by Schapire, consists of the combination of error correcting output codes with boosting techniques [4]; this ensemble method shows good performances on different benchmark machine learning problem [19].

From a statistical standpoint ECOC methods can be viewed as an approximation of a Bayes classifier: James

demonstrated that asymptotically, as the number of dichotomisers approaches infinity, the ECOC classifier will become Bayes consistent (i.e. it always classifies to the Bayes class when the base learner is the Bayes classifier), provided that a random coding matrix is used [20]. In the same perspective Berger showed that randomly selected decomposition matrices are likely to have pairwise well-separated codewords, that is high error recovering capabilities [21]. Variants of the original ECOC algorithm have been proposed, as circular ECOC [22] to reduce the sensitivity to codeword selection, or binary labelling techniques [23] to reduce the correlation between the base learners.

The good generalisation properties of ECOC methods have been explained through the reduction of both bias and variance [21,15] and by interpreting them as large margin classifiers [24,6]. Application of ECOC methods in several domains have shown improvements over standard k -way classification methods. For instance, ECOC was successfully applied to classify cloud types [25], for text classification [21,26], for text-to-speech synthesis [27], to classify olive oils by means of electric noses [28], for face verification [29], and to classify malignant and normal tissue using gene expression DNA microarray data [30].

ECOC methods present several open problems concerning their properties and the factors affecting their effectiveness. The trade-off between error recovering capabilities and learnability of the dichotomies induced by the decomposition scheme have been tackled in several works [24,31], but an experimental evaluation of the trade-off has to be performed to achieve a better understanding of this phenomenon.

A related problem is the analysis of the relationship between codeword length and performances: some preliminary results seem to show that long codewords improve performance [26]. Another problem, not sufficiently investigated in literature [21,26,32], is the proper selection of dichotomic learning machines for the decomposition unit.

Designing codes for a given multiclass problem is another interesting open problem. A greedy approach is proposed in Mayoraz and Moreira [33], and a method based on soft weight sharing to learn error correcting codes from data is presented in Alpaydin and Mayoraz [34]. In Crammer and Singer [35] it has been shown that given a set of dichotomisers the problem of finding an optimal decomposition matrix in P-complete: by introducing continuous codes and casting the design problem of continuous codes as a constrained optimisation problem, we can achieve an optimal continuous decomposition using standard optimisation methods.

In this paper, we tackle some of the open problems concerning ECOC methods, and we try to experimentally analyse the factors affecting the effectiveness of ECOC classifiers. In particular, we study if the architecture of ECOC learning machine influences the dependence among codeword bit errors and the performance of the overall multiclassifier. Moreover we experimentally compare different

decoding functions and different base learners in order to evaluate their influence on the generalisation error of the ECOC ensembles. The relationships between ensemble accuracy, base learner accuracy and the error correction power of ECOC codes are experimentally analysed to understand if the error recovering capabilities of ECOC codes can by itself explain the good generalization capabilities of ECOC methods, or if they are the result of complex interactions between the error recovering power of ECOC, the complexity of the induced dichotomies and the accuracy of the base learners composing the ensemble.

In the next session an overview of the ECOC methods is given. Then we experimentally analyse the factors affecting the effectiveness of ECOC methods. A discussion on the results and an outline of future developments of this work concludes the paper.

ECOC for multiclass learning problems

In this section we outline the main characteristics of ECOC methods for multiclass classification. ECOC methods code classes through binary strings and exploit the redundancy of the resulting coding schemes to reduce the classification error. They are characterised by a decomposition of a multiclass problem in a set of dichotomic problems any by a successive reconstruction of the original multiclass problem. Two main architectures are feasible, one based on a single learning machine, and another one based on an ensemble of classifiers.

Decomposition of a multiclass classification problem

In a classification problem based on decomposition methods [36], usually we code classes through binary strings, or codewords. A coding process is a mapping $\mathcal{M}: \{C_1, \dots, C_k\} \rightarrow \{s_1, \dots, s_k\}$ from the set of classes to the set of binary strings. Each string s_i , $1 \leq i \leq k$ must univocally determine its corresponding class.

Let be $\mathcal{P}: \mathbf{X} \rightarrow \{C_1, \dots, C_k\}$ a K classes polychotomy (or K -polychotomy), where \mathbf{X} is the multidimensional space of the features and C_1, \dots, C_k are the labels of classes. The decomposition of the K -polychotomy generates a set of L dichotomies f_1, \dots, f_L . Each dichotomy f_i subdivides the input patterns in two complementary superclasses C_i^+ and C_i^- , each of them grouping one or more classes of the K -polychotomy. Let be also $D = [d_{ik}]$ a *decomposition matrix* of dimension $L \times K$ that represents the decomposition, connecting classes C_1, \dots, C_k to the superclasses C_i^+ and C_i^- identified by each dichotomy. An element of D is defined as:

$$d_{ik} = \begin{cases} +1 & \text{if } C_k \subseteq C_i^+ \\ -1 & \text{if } C_k \subseteq C_i^- \end{cases} \quad (1)$$

A learning algorithm produces an hypothesis $\hat{\mathbf{f}}(x) = [\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_L(x)]$ of the unknown functional $f(x) = [f_1(x), f_2(x), \dots, f_L(x)]$, where $[f_1(x), f_2(x), \dots, f_L(x)]$

is the codeword s associated by the coding function \mathcal{M} to the class of the input pattern x . When a polychotomy is decomposed into dichotomies, the learning task for $\mathbf{f}(x)$ is reduced to the learning of each $f_i: \mathbf{X} \rightarrow \{-1, 1\}$ through the set of the dichotomisers: it consists in labelling some classes with +1 and others with -1. Each dichotomiser is trained to learn f_i , that associates patterns belonging to class C_k with values d_{ik} of the decomposition matrix D , producing the hypothesis \hat{f}_i .

In the decomposition matrix, rows correspond to dichotomisers tasks and columns to classes. In this way, each class is univocally determined by its specific codeword. For instance, considering a decomposition matrix for a four class classification problem with 7-bit class coding (Table 1), the task of the second dichotomiser, namely f_2 , consists in separating the patterns belonging to classes C_1 and C_4 from the patterns of class C_2 and C_3 . The third column of the decomposition matrix represents the codeword $[-1, -1, +1, +1, -1, -1, +1]$ associated to the class C_3 .

ECOC decomposition tries to maximise error recovering capabilities through the maximisation of the minimum distance between each couple of codewords [15,36].

Several methods for generating ECOC codes have been proposed: exhaustive codes, randomised hill climbing [12], Hadamard and BCH codes [16,37], and random codes [20], but open problems are still the joint maximisation of distances between rows and columns of the decomposition matrix.

Reconstruction and decoding

After the training of the dichotomisers \hat{f}_i , their outputs are used to reconstruct the polychotomy to determine the class $C_i \in \{C_1, \dots, C_k\}$ of the input patterns, using a suitable measure of similarity.

Learning machines constructed by ECOC are made up by a *Decomposition Unit* and a *Decision Unit*. The Decomposition Unit analyses the input patterns and calculates the codewords using an assigned decomposition scheme generated by a suitable algorithm. This unit computes:

$$\hat{\mathbf{f}}(x) = [\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_L(x)] \quad (2)$$

Table 1 ECOC decomposition matrix example

Dichotomisers tasks	Columns: class codewords			
	C_1	C_2	C_3	C_4
f_1	+1	-1	-1	-1
f_2	+1	-1	-1	+1
f_3	+1	-1	+1	-1
f_4	+1	-1	+1	+1
f_5	+1	+1	-1	-1
f_6	+1	+1	-1	+1
f_7	+1	+1	+1	-1

The Decision Unit decodes the computed codeword $\hat{s} = [\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_L(x)]$, mapping it to the associated class. This unit computes the function $\mathcal{D}: S \rightarrow C$:

$$\mathcal{D}(\hat{s}) = \mathcal{D}[\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_L(x)] \quad (3)$$

where S is the set of the computed codewords, C is the set of the classes, $\hat{f}_i(x)$ are the hypotheses returned by the learning algorithm, and \mathcal{D} is a suitable decoding function.

The Decision Unit decodes the codeword computed by the Decomposition Unit, choosing the class whose codeword is more similar to that computed by the set of dichotomisers. So, the decoding function $\mathcal{D}(s)$ can be implemented by a maximization of a similarity measure between the computed s codeword and the effective codewords s_i , $1 \leq i \leq K$ associated to the classes:

$$class_{out} = \mathcal{D}(\hat{s}) = \arg \max_{1 \leq i \leq K} Sim(\hat{s}, s_i) \quad (4)$$

where $class_{out}$ is the class computed by the polychotomiser, s_i is the codeword of class C_i , the vector \hat{s} is the codeword computed by the set of dichotomisers, and $Sim(x, y)$ is a general similarity measure between two vectors x and y . This similarity measure can be a Hamming distance if the outputs of the dichotomisers $\hat{s} = [\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_L(x)]$ are discrete or the inner product, or L_1 or L_2 norm distances for dichotomisers with continuous outputs. For instance, using the inner product we have:

$$\mathcal{D}_{prod}(\hat{s}) = \arg \max_{1 \leq i \leq K} (\hat{s}, s_i) \quad (5)$$

Dietterich and Bakiri [16,17] proposed the Error Correcting Output Coding (ECOC) decomposition scheme with the aim of improving the generalisation capabilities of NETtalk classifier systems based on distributed output codes [18]: Coding the classes by codewords suggests the idea of adding *error recovering* capabilities to decomposition methods to obtain classifiers less sensitive to noise [15,32]. This goal is achieved by means of the redundancy of the coding scheme, as shown by coding theory [38].

The error-recovering capabilities of ECOC codes depend mainly on column separation, i.e. the distance between codewords must be increased, according to an assigned measure. The maximal number of error Max_{err} that can be corrected in an ECOC based decomposition is [15]:

$$Max_{err} = \lfloor \frac{\Delta_D - 1}{2} \rfloor \quad (6)$$

where Δ_D is the minimal Hamming distance between each pair of columns of the decomposition matrix D .

Design of ECOC classifiers

There are two main approaches to the design of a classifier using ECOC codes, based on the features of the Decomposition Unit. The first approach, that we call *Monolithic Classifiers*, makes use of a Decomposition

Unit based on a multiple inputs multiple outputs (MIMO) learning machine exploiting the decomposition in an implicit way. *Monolithic Classifiers* are, for example, MIMO MLPs or MIMO decision trees [11, 12].

The second approach, that we call *Parallel Classifiers*, makes use of a Decomposition Unit leading to the distribution of the learning task among separated and independent dichotomisers that can be implemented through, e.g. Support Vector Machines [39], multiple inputs single output (MISO) MLP or dichotomic decision trees. We call the resulting learning machines *Parallel Linear Dichotomisers (PLD)* if the dichotomisers used for implementing the dichotomisers are linear (as in Alpaydin and Mayorez [40]), or *Parallel Non-linear Dichotomisers (PND)* if the dichotomisers are non-linear [12, 41].

Parallel Non-linear Dichotomisers (PND) are multi-classifiers based on the decomposition of polychotomies into dichotomies, using dichotomisers solving their classification tasks independently from each other [41]. In the decomposition unit each dichotomiser is implemented by a separate *non-linear* learning machine, and learns a different and specific dichotomic task using a training set common to all the dichotomisers. The decision unit can use a L_1 norm or another similarity measure between codewords to predict classes of unlabeled patterns. *PND* have been implemented with decision trees [12,15,19] or dichotomic MLP [41].

Parallel Linear Dichotomisers (PLD) are also multi-classifiers based on decomposition of polychotomies into dichotomies, but each dichotomiser is implemented by a separate linear learning machine [40].

It is worth noting that classifiers based on decomposition methods and classifiers based on ensemble averaging methods [2,3] or bagging and boosting [7, 4] share the idea of using a set of learning machines acting on the same input and recombining their outputs in order to make decisions; the main difference lies in the fact that in classifiers based on decomposition methods the task of each learning machine is specific and different from that of the others.

Effectiveness of ECOC methods

In this section we experimentally analyse the factors affecting the effectiveness of ECOC methods. In particular we focus on the following items:

1. Architecture of the decomposition unit.
2. Dependency among codeword bits coding the classes.
3. Decoding function selected for the decision unit.
4. Relationships between ensemble accuracy, base learner accuracy and error correcting power.

In the following sections we address each problem separately.

Architectures

In this section we study two different architectures for ECOC learning machines, considering decomposition units composed by a single monolithic MLP, that learns the codewords as a whole, and by an ensemble of dichotomic MLP, each learning a different bit of the codewords coding the classes.

We apply also the widely used One-Per-Class¹ (OPC) [42, 43] decomposition scheme as reference comparison, both using MLP monolithic and MLP parallel classifiers decomposition units.

PND are implemented by a set of multi-layer perceptrons with a single hidden layer, acting as dichotomisers, and *PLD* are implemented by a set of single layer perceptrons.

Monolithic MLP are built using a single hidden layer and sigmoidal activation functions, both in hidden and output neurons. The number of neurons of the hidden layer amounts from ten to one hundred according to the complexity of the data set to be learned. The base learners of the *PND* and the monolithic MLP have been trained using the backpropagation algorithm with fixed learning rate.

We have compared the performances of the three ECOC learning machines using different data sets, both real and synthetic, as shown in Table 2. The data set *p6* and *p9* (available by anonymous ftp at ftp://ftp.disi.unige.it/person/ValentiniG/Data), are synthetic and composed by normal distributed clusters of data. The set *p6* contains six classes with no overlapping regions, while the regions of the nine classes of *p9* hardly overlap. *Glass*, *letter* and *optdigits* data sets are from the UCI repository [44].

In the experimentation we have used exhaustive [12] and BCH ECOC generation algorithms [16]. ECOC exhaustive algorithms select among all possible 2^K dichotomies the $2^{K-1} - 1$ ones that are not equivalent and not trivial². ECOC obtained by exhaustive algorithms are Bayes consistent [20], i.e. if the component dichotomisers approximate the Bayes optimal discriminant function, then the overall polychotomiser will produce an optimal

Table 2 Data sets general features. The *glass*, *letter* and *optdigits* data sets are from the UCI repository [44]

Data set	Number of attributes	Number of classes	Number of training samples	Number of testing samples
<i>p6</i>	3	6	1200	1200
<i>p9</i>	5	9	1800	5-fold cross-val
<i>glass</i>	9	6	214	10-fold cross-val
<i>letter</i>	16	26	16000	4000
<i>optdigits</i>	64	10	3823	1797

¹ In a One-Per-Class (OPC) decomposition scheme, each dichotomiser f_i have to separate a single class from all the others.

² Dichotomies f' and f such that $f' = f$ or $f' = -f$ are equivalent, while for trivial dichotomies $f^{-1}(+1) = \emptyset$ or $f^{-1}(-1) = \emptyset$ holds.

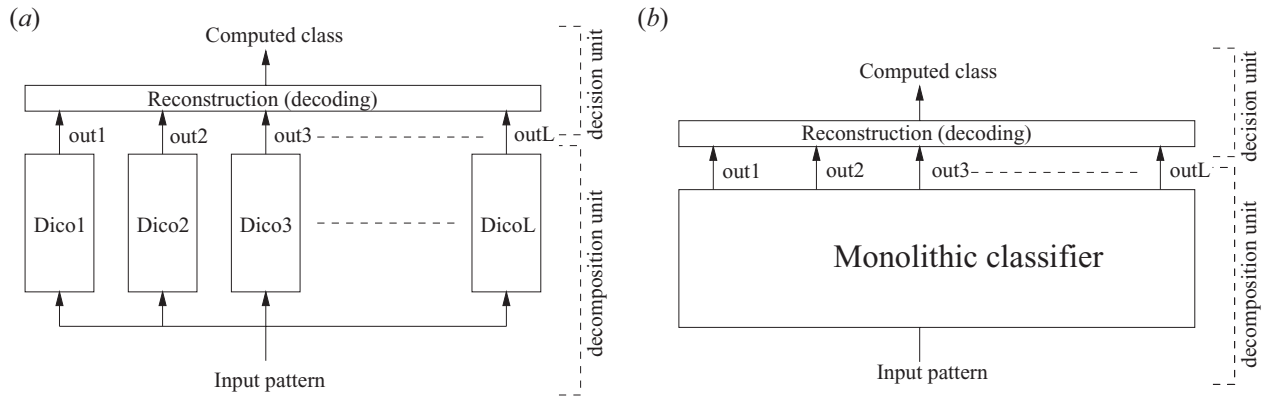


Fig. 1 Design of output coding learning machines: monolithic (a) and parallel ensemble (b)

Bayes classification. The shortcoming of the algorithm is the exponential growth of the codeword lengths with the cardinality of the classes. BCH ECOC are computed using an algebraic method based on a polynomial representation of finite Galois fields [16]. BCH ECOC are not Bayes optimal, but allow to generate ECOC codewords of tractable length. We have modified the original algorithm removing duplicate rows of the decomposition matrix or trivial dichotomies, and deleting rows with a Hamming distance equal or less than an assigned threshold. This modified version of the algorithm has produced 7 bits ECOC codes for the data sets *p6* and *glass*, 15 bits codes for *p9* and *optdigits* and 30 bit ECOC codes for *letter* data set.

The programs used in our experiments have been developed using *NEUROjects* [45], a C++ library for neural networks development. In the experimentation we used resampling methods, using a single pair of training and testing data set or *k-fold cross validation*, to estimate the generalisation error of ECOC monolithic and ECOC *PND* ensembles.

Figure 2, shows the performance of MLP, PLD and *PND* over the considered data sets. ECOC MLP monolithic classifiers do not outperform standard MLP (Fig. 2(a)). In Dietterich and Bakiri [12] similar results have been obtained over the same data set *letter* we used.

Concerning *PLD* (Fig. 2(b)), over data sets *p6*, *p9*, and *optdigits* there is no significant statistical difference among OPC and ECOC decomposition, while over *glass* *PLD* ECOC outperforms all other types of polychotomisers, but with *letter* *PLD* OPC achieve better results.

Considering *PND* (Fig. 2(c)), for data sets *p6* and *optdigits* no significant differences among OPC and ECOC *PND* can be noticed. Over the *p9* data set, ECOC shows expected errors significantly smaller than OPC. Expected errors over *glass* and *letter* data sets are significantly smaller for ECOC compared with OPC. So we can see that, ECOC *PND* show expected error rates significantly lower than OPC *PND*.

The differences in performances between OPC and ECOC decomposition schemes are considered statistically significant if their confidence level is less or equal to 0.05

according to *McNemar* test [46] or *k fold cross validated paired t test* [47].

PLD show remarkable higher errors over all data sets, and in particular they fail over *p9*. Summarising, the expected errors are significantly smaller for *PND* compared with direct monolithic MLP classifiers and *PLD*, and ECOC outperforms OPC decomposition only in *PND* ensembles.

Dependency among codeword bits

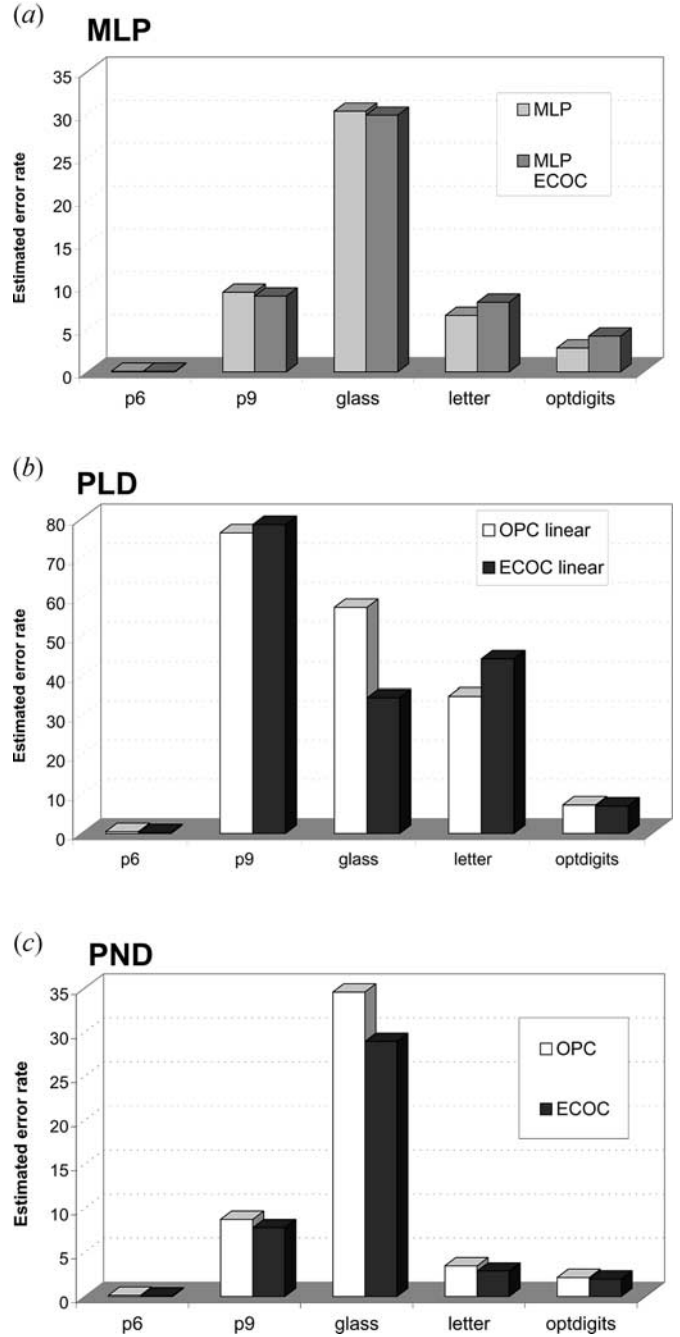
Several authors noted that the dependence among output errors affect the effectiveness of ECOC methods [15,32,48], and recently we provided a quantitative evidence of this fact [49]. In the domain of the serial transmission of messages coded as sequences of bits, Peterson and Weldon [37] showed that if errors on different codeword bits are dependent, the effectiveness of error correcting code is reduced. Transferring these outcomes in the framework of classification problems, if a decomposition matrix contains very similar rows (dichotomies), each error of an assigned dichotomiser will be likely to appear in the most correlated dichotomisers, thus reducing the effectiveness of ECOC.

To quantitatively evaluating the dependence among output errors of the decomposition unit of ECOC learning machines, we used mutual information based measures proposed in Masulli and Velentini [50].

Mutual information, being a special case of the Kullback-Leibler divergence between two distributions, measures the matching between the joint probability density distribution and the product of the marginal probability density distribution of the output errors. If we have a complete matching, the mutual information is 0 and the output errors are independent, otherwise the higher the value of the mutual information between output errors is, the higher the dependence between them will be.

Figure 3 shows the compared *mutual information error index* Φ_R [50] between *monolithic* and *PND* ECOC learning machines considering 4 different data sets. It is defined as

Fig. 2 Performance of ECOC and OPC MLP, PLD and PND. (a) Monolithic MLP (b) PLD (c) PND



$$\Phi_R = \sum_{i=1}^l \sum_{j=1}^l I_E(e_i, e_j) \quad (7)$$

where $I_E(e_i, e_j)$ is the mutual information between the errors of the i th and j th output of the decomposition unit.

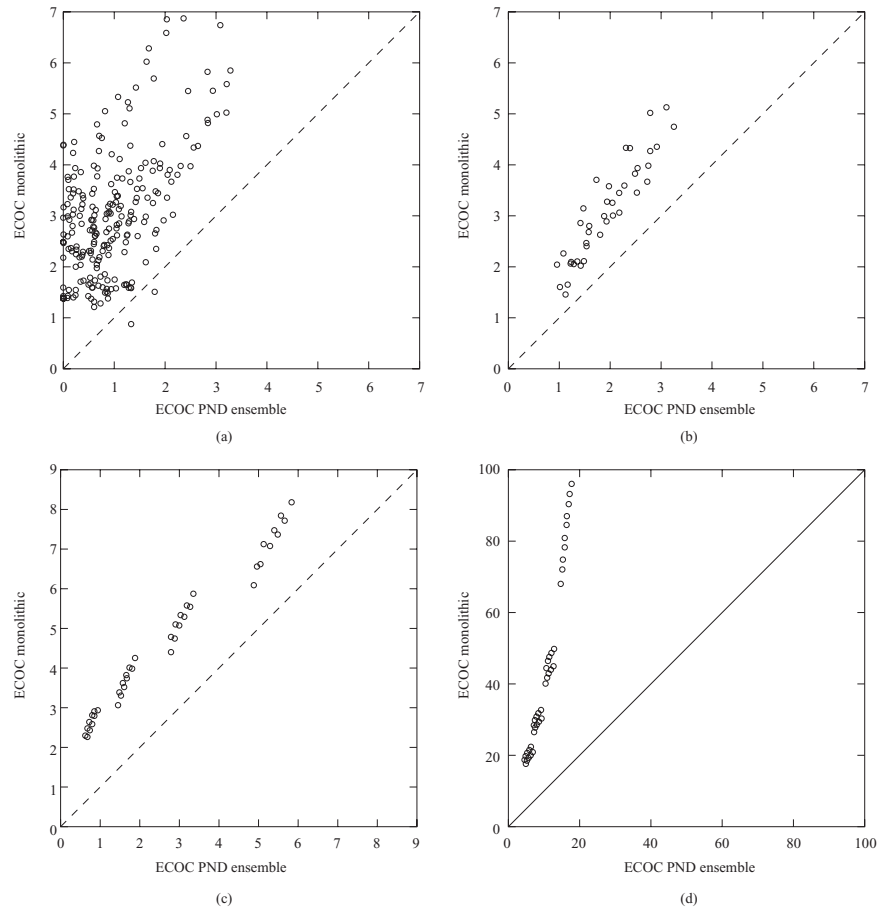
On the axes are represented, Φ_R values of ECOC *monolithic* and ECOC *PND* learning machines. This index measures the sum of the mutual information of the output errors between all the output pairs of the learning machines, giving a computable quantity to estimate the dependence between codeword bit errors: a high value of Φ_R corresponds to an high dependence between output errors, and vice versa. Each point corresponds to ECOC

learning machines implemented with MLP with different number of hidden units and using different partitions of the output error. On all the data sets about all the points are above the dotted line, i.e. all the values of Φ_R are greater for ECOC *monolithic* compared with ECOC *PND*. The results show that *monolithic* architectures are affected by a higher dependence among codeword bit errors, confirming the results obtained in Masulli and Valentini [49].

Relationships between ensemble accuracy and decoding function

In this section we experimentally study if the choice of a particular decoding function affects the performance of

Fig. 3 Compared mutual information specific error matrix indices Φ_R between ECOC *monolithic* and *PND* learning machines on d5 (a), glass (b), optdigits (c) and latter (d) data sets. d5 is a synthetic data set available at <http://ftp.disi.unige.it/person/ValentiniG/Data/d5.tgz>



ECOC MLP ensembles. We analyse also the effectiveness of the ECOC ensemble varying the minimum Hamming distance with fixed length codewords, as the error recovering capabilities of the ECOC ensembles depend critically on the minimum Hamming distance between code-words (Eq. (6)).

The decoding in Output Coding methods is performed using similarity measures between the computed codeword and the codeword coding the classes (Sect. 2). We consider three commonly used decoding functions based on Hamming distance and L_1 and L_2 norm distance.

Given a $n \times k$ decomposition matrix D (Eq. (1)), we indicate with D_{ij} the i th bit of the j th codeword coding the classes, and with $C_i(x)$ the output computed by the i th dichotomiser on the input $x \in \mathbb{R}^d$. If $C_i(x) \in \{-1, +1\}$, then the decoding function based on the Hamming distance is:

$$\mathcal{D}_{Hamming}(x) = \arg \min_j \sum_{i=1}^n \frac{1}{2} |D_{ij} - C_i(x)| \quad (8)$$

If $C_i(x) \in \mathbb{R}^d$, then the decoding function based on the L_1 norm distance is:

$$\mathcal{D}_{L_1}(x) = \arg \min_j \sum_{i=1}^n |D_{ij} - C_i(x)| \quad (9)$$

and the decoding function based on the L_2 norm distance is:

$$\mathcal{D}_{L_2}(x) = \arg \min_j \sum_{i=1}^n (D_{ij} - C_i(x))^2 \quad (10)$$

We generated ECOC decomposition matrices with constrained random algorithms. The random generation have been constrained in order to eliminate trivial dichotomies (e.g. rows with all +1 or all -1), and equal or complementary rows to assure the absence of equal and equivalent symmetric dichotomies, and to achieve a desired minimum Hamming distance between the columns (codewords) of the decomposition matrix³.

The generalisation error of the ECOC ensembles have been estimated using 5-fold cross validation. We merged the training and test sets of the optdigits and image-segmentation data sets from the UCI repository to perform a five-fold cross validation on the overall merged data set. We used also p20, a synthetic data set generated through the *NEUROObjects* application dodata. It is composed by 20 3-dimensional classes, and each class is characterised

³ The C++ classes implementing the ECOC random algorithms are available at <http://ftp.disi.unige.it/person/ValentiniG/SW/ECOC/random>

by three disjoint clusters of data normally distributed with diagonal covariance matrices⁴.

The comparison of the estimated generalisation error of the ECOC ensembles shows that the decoding functions based on the L_1 and L_2 norm distances outperform the decoding based on the Hamming distance (Figs 4–6). Only on the *p20* data set with 50-bit ECOC ensembles with linear perceptrons as base learners there is no difference between L_1 norm and Hamming distance based decoding, but in this task the ECOC ensemble clearly fails, performing a sort of random guessing. The same figures show also that there is no significant difference between L_1 and L_2 norm distance decoding, considering the proposed data sets.

Considering the relationship between the estimated generalisation error and the minimum Hamming distance (MHD) between the codewords, at first glance, we could expect a monotonic decrement with the MHD. But the trends seem to be more complex: they are relatively irregular (Figs 4–6), and only in a few cases (Figs 4(c),

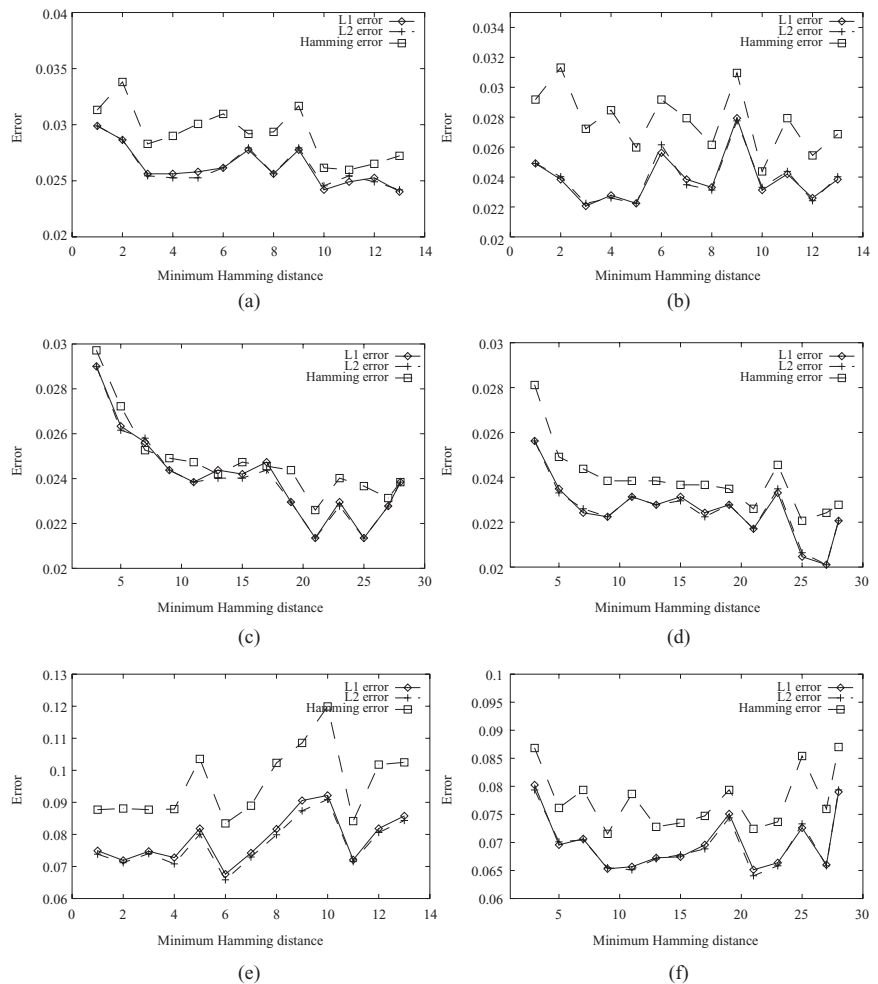
(d)) we can observe a monotonic decrement of the error with MHD.

Relationships between ensemble and base learner accuracy and error correcting power

The previous section showed that the increment of the MHD is not by itself sufficient to reduce the error of the ensembles. In this section we study how the characteristics of the base learner, and in particular its accuracy affects the overall performance of the ECOC ensemble.

As expected, the selection of the base learner affects the overall performance of the ensemble (Fig. 7). Indeed, using MLP base learner with different number of hidden units greatly influences the estimated generalisation error of the ensemble. For instance, on the *p20* data set, ECOC ensembles having as base learners MLP with 24 hidden units are able to halve the estimated generalisation error

Fig. 4 Relationships between ECOC minimum Hamming distance among codewords and ensemble errors using decoding functions based on L_1 and L_2 norm and Hamming distance with the *optdigits* data set. (a) 32 bit ECOC codewords using MLP base learners with 4 hidden units (b) 32 bit ECOC using MLP with 10 hidden units (c) 64 bit ECOC using MLP with 3 hidden units (d) 64 bit ECOC using MLP with 10 hidden units (e) 32 bit ECOC using linear perceptrons base learners (f) 64 bit ECOC using linear perceptrons base learners



⁴ *p20* is available at <http://ftp.disi.unige.it/person/ValentiniG/Data/p20.tgz>

Fig. 5 Relationships between ECOC minimum Hamming distance among codewords and ensemble errors using decoding functions based on L_1 and L_2 norm and Hamming distance with the *image-segmentation* data set. (a) 32 bit ECOC codewords using MLP base learners with 3 hidden units (b) 32 bit ECOC codewords using MLP base learners with 10 hidden units (c) 32 bit ECOC codewords using MLP base learners with 25 hidden units (d) 32 bit ECOC codewords using linear perceptron base learners

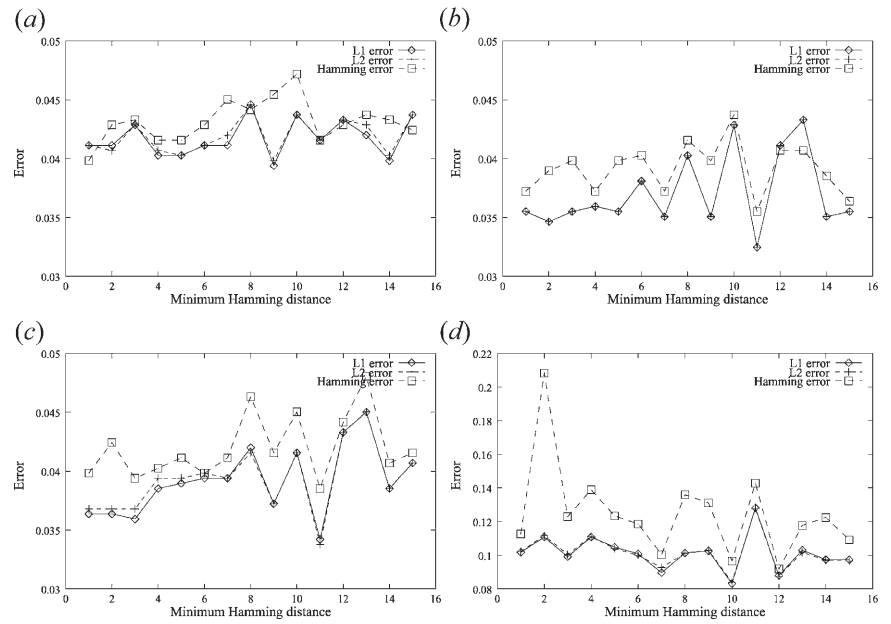
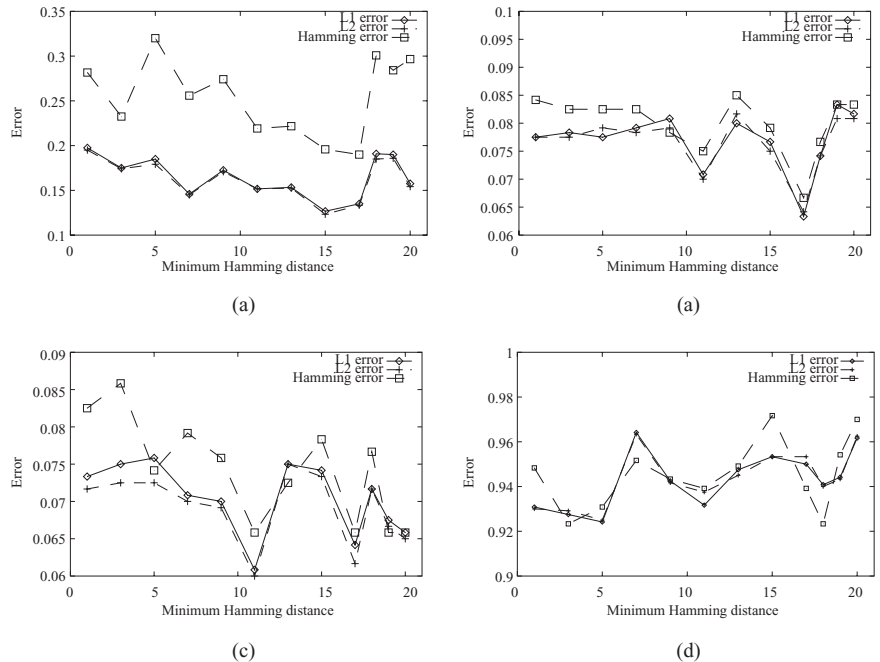


Fig. 6 Relationship between ECOC minimum Hamming distance among codewords and ensemble errors using decoding functions based on L_1 and L_2 norm and Hamming distance with the *p20* data set. (a) 50 bit ECOC codewords using MLP base learners with 6 hidden units (b) 50 bit ECOC codewords using MLP base learners with 12 hidden units (c) 50 bit ECOC codewords using MLP base learners with 24 hidden units (d) 50 bit ECOC codewords using linear perceptron base learners



of ECOC ensembles with MLP with 6 hidden units (Fig. 7(d)). Also on the *optdigits* and *image-segmentation* data set (Figs 7(a)–(c)) the selection of different MLP as base learner has a significant impact on the performance of the ensemble.

In order to investigate why the selection of different base learners affects in a so significant manner the performance of the ensemble, we analysed the relationship between the overall ensemble error, the average base learner error and the minimum Hamming distance between the codewords.

With the *optdigits* data set the average base learner error shows an irregular trend, with minima of the error increments with the minimum Hamming distance

(Fig. 8), while the ensemble error tends to decrease with the MHD, especially using long codewords (Figs 8(c), (d), (f)). More irregular patterns are observed in the *image-segmentation* data set: the average base learner error increases with the Hamming distance, but in an irregular way and the ensemble error rate oscillates around 0.04 and 0.035 (Fig. 9). With the *p20* synthetic data set the base learner average error initially grows, then goes down and hence increments with the MHD, while the ensemble error decreases only using base learners with 6 hidden units: using more complex base learners the ensemble error decreases for 11 and 17 bit MHD (Fig. 10). Note that the ensemble

Fig. 7 Relationships between ECOC minimum Hamming distance and ensemble error using different base learners. (a) *optdigits* data set (32 bit codewords) (b) *optdigits* data set (64 bit codewords) (c) image-segmentation data set (d) p20 data set

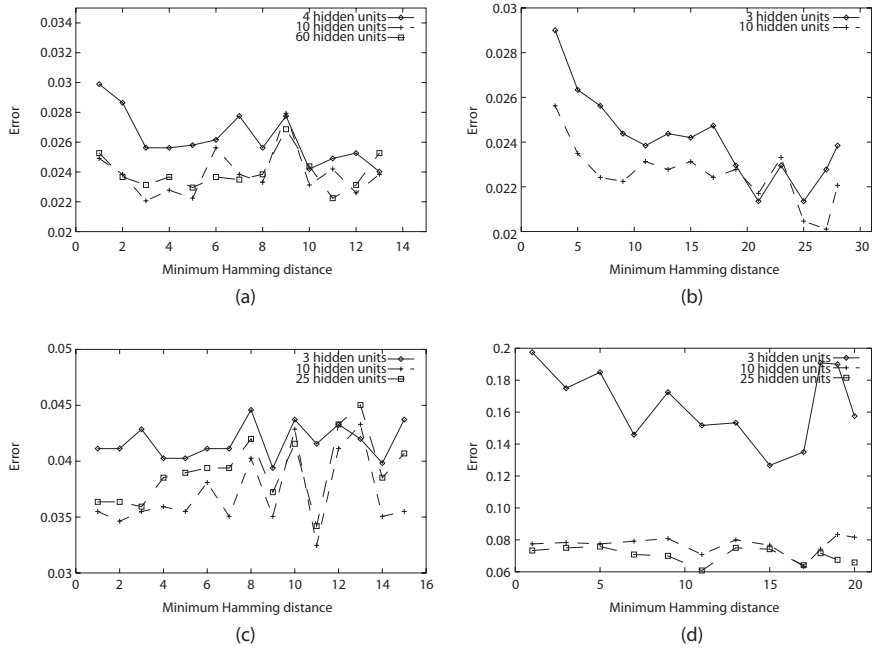


Fig. 8 Relationships between ECOC minimum Hamming distance among codewords, ensemble error and average base learner error with the *optdigits* data set. (a) 32 bit ECOC codewords using MLP base learners with 4 hidden units (b) 32 bit ECOC using MLP with 10 hidden units (c) 64 bit ECOC using MLP with 3 hidden units (d) 64 bit ECOC using MLP with 10 hidden units (e) 32 bit ECOC using linear perceptrons base learners (f) 64 bit ECOC using linear perceptrons base learners

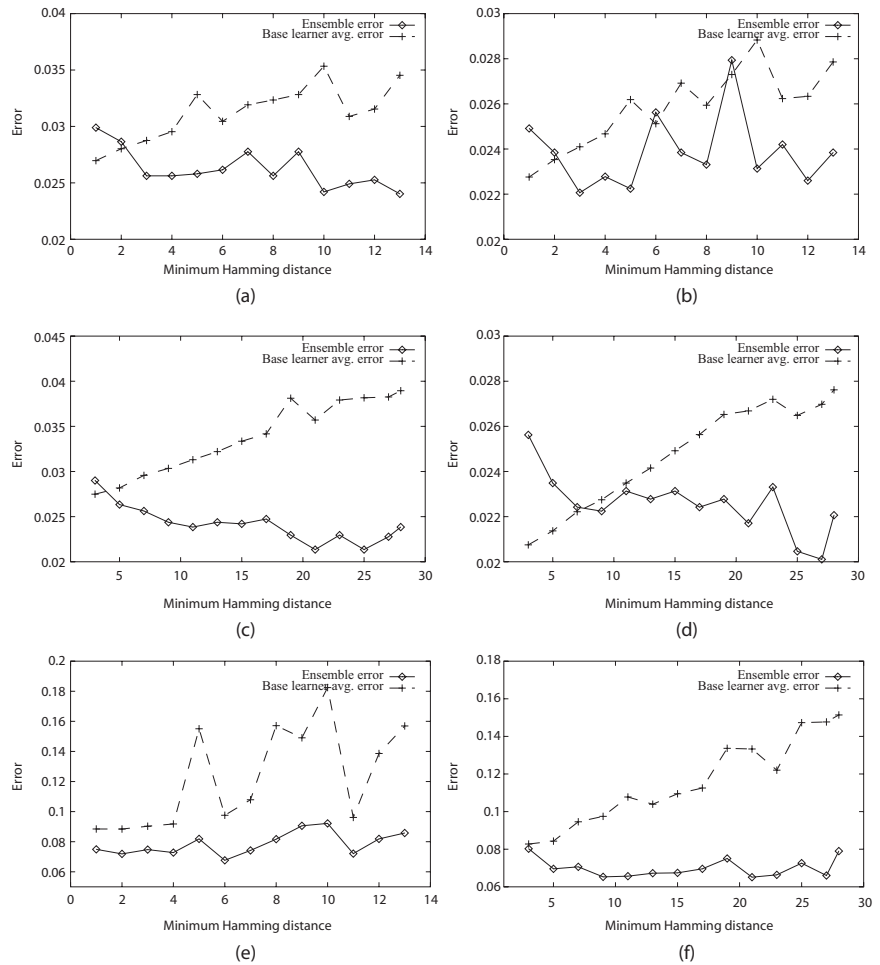


Fig. 9 Relationships between ECOC minimum Hamming distance between codewords, ensemble error and average base learner error with the *image-segmentation* data set. (a) 32 bit ECOC codewords using MLP base learners with 3 hidden units (b) 32 bit ECOC codewords using MLP base learners with 10 hidden units (c) 32 bit ECOC codewords using MLP base learners with 25 hidden units (d) 32 bit ECOC codewords using linear perceptron base learners

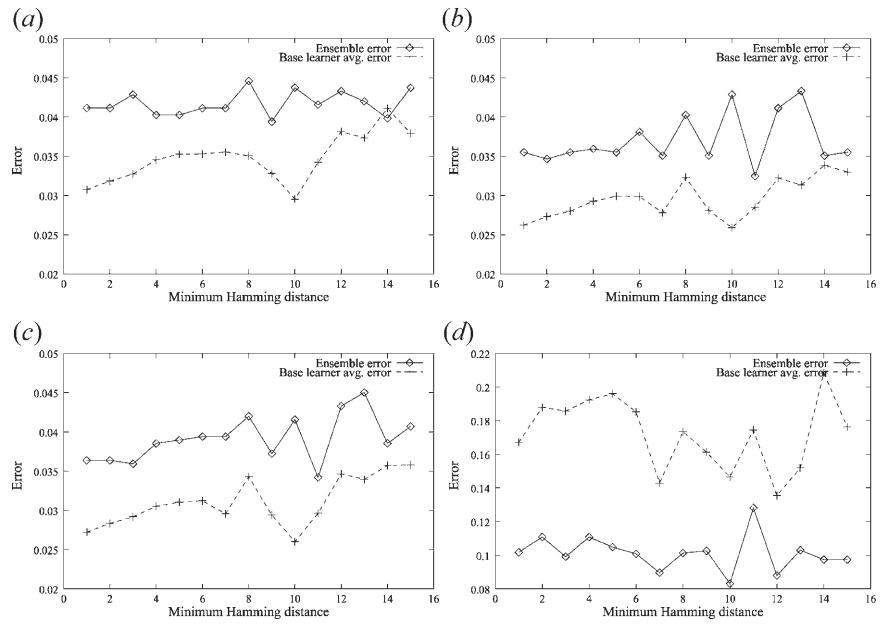
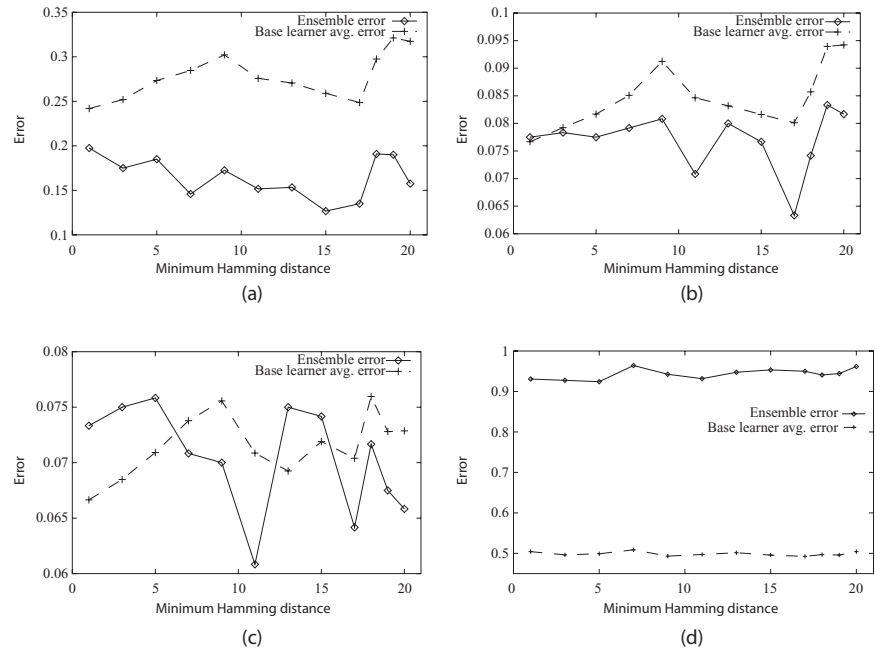


Fig. 10 Relationships between ECOC minimum Hamming distance between codewords, ensemble error and average base learner error with the *p20* data set. (a) 50 bit ECOC codewords using MLP base learners with 6 hidden units (b) 50 bit ECOC codewords using MLP base learners with 12 hidden units (c) 50 bit ECOC codewords using MLP base learners with 24 hidden units (d) 50 bit ECOC codewords using linear perceptron base learners



error partially follows the average base learner error, with a decrement for increasing MHD values, but sometimes we can also observe opposite trends of the average base learner and ensemble errors (Figs 8–10).

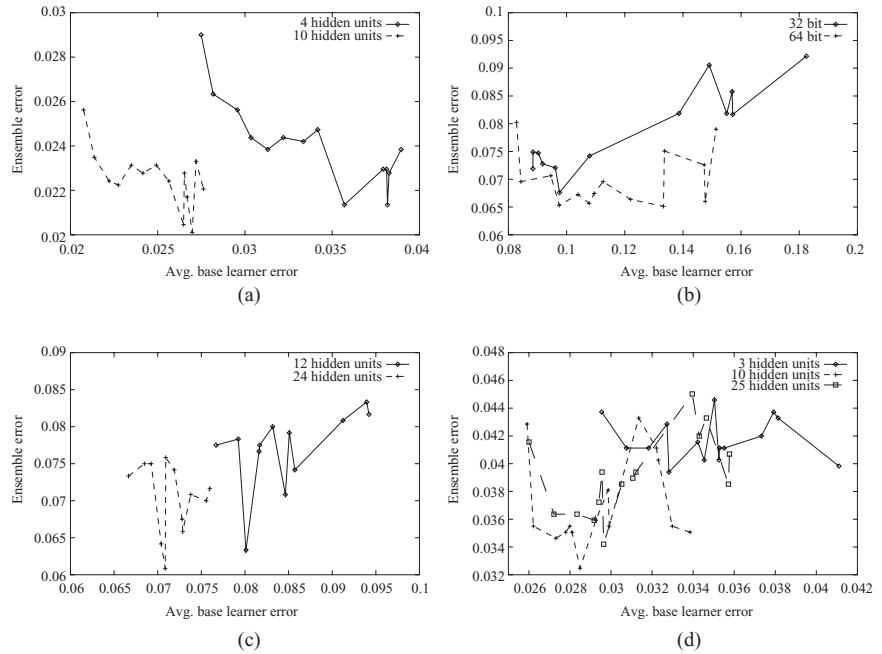
The graphs of Fig. 11 represent directly the ensemble error in function of the average base learner error. Each point in the graph corresponds to a different MD between codewords. Considering different data sets and different base learners we observe very different trends: with the *opt digits* data set and MLP base learners we have a decrement of the ensemble error even if the average base learner error increases (Fig. 11(a)), but with the same data set using linear perceptrons base learners we have an increment of the error, especially with 32-bit codewords

(Fig. 11(b)). With the *p20* and *image-segmentation* data sets the trends seem to be more complex, without a clear relationship between ensemble and average base learner accuracy (Figs 11(c), (d)). Summarising, we cannot observe simple relationships between the ensemble error and the average base learner error with respect to the MHD between codewords.

Discussion

The experimental results of the previous section show that a set of different factors affect the effectiveness of the ECOC methods. The architecture of the ECOC learning

Fig. 11 Relationships between Average base learner error and ensemble error. Each point represents a different minimum Hamming distance between codewords. (a) optdigits: 64 bit ECOC codewords using MLP base learners (b) optdigits: 64 bit ECOC using linear perceptron base learners (c) p20: 50 bit ECOC using MLP base learners (d) image-segmentation: 32 bit ECOC using MLP base learners



machine, the dependency among codeword bits, the type of the decoding function, the error recovering power of the decomposition scheme, and the accuracy of the dichotomiser interact between them and contribute to the effectiveness of the ECOC methods.

The design of the decomposition unit of the ECOC learning machine affects the performances, as shown by our experimentation. It has been stated [11,15] that ECOC classifiers should be preferred to OPC classifiers, as they reduce error bias and variance more than standard classifiers; their experimental results confirm these hypotheses, with the exception of some cases over complex data sets (such as *letter* from UCI repository) where OPC MLP classifiers perform better than ECOC MLP. Our experimentation has pointed out that not always ECOC monolithic MLP outperform OPC MLP classifiers, while we have found a significant difference between ECOC and OPC *PND* performances (Fig. 2). We suppose that the better performances of *PND* can be explained considering on one hand that their dichotomisers are less complex than ECOC *monolithic* learning machines, achieving by this way better generalisation capabilities. On the other hand, ECOC *monolithic* learning machines introduce more correlation among codeword bits. In fact in *PND* each codeword bit is learned and computed by its own MLP, specialised for its particular dichotomy, while in monolithic classifiers each codeword bit is learned and computed by a linear combination of hidden layer outputs pertaining to one and only shared multi-layer perceptron. Concerning PLD, we point out that the error recovering capabilities induced by ECOC are counter-balanced by higher error rates of linear dichotomisers (Figs 5(d) and 6(d)). Hence, interdependence among monolithic MLP ECOC outputs lowers the effectiveness of ECOC codes for this kind of classifiers. This is conformed also by the quantitative evaluation of the dependence among code-

word bit errors we performed using mutual information based measures: ECOC *PND* ensembles show a lower dependence compared with monolithic ECOC MLP.

The results showed also that the decoding function plays an important role: indeed L_1 and L_2 norm distance seem to be well-suited for the decoding, while the Hamming distance based decoding function achieves worse results. This is not surprising, as L_1 and L_2 norms exploit the “confidence” in the prevision of each base learner, while Hamming decoding discards all the information except the hard membership to a class.

The choice of proper dichotomisers, well-suited for a given decomposition, affects also the performances of ECOC ensembles, as shown by our experimental results (Fig. 7): in general complex classification problems need more complex dichotomisers, but overfitting phenomena can also arise.

The error recovering power of ECOC methods depends on the MHD between codewords (if the output errors of the decomposition unit are independent). Our experiments show that if we use fixed length codewords an increment of the MHD does not lead necessarily to improved performances of the ECOC ensemble. This can be explained considering that different codewords induce different dichotomies. The dichotomies can or cannot be hard learnable depending on the structure of the data and on the type of the base learner used. The learn ability is partially reflected by the average base learner error. As shown by our experiments there is not a simple relationship between ensemble error, average base learner error and the MHD. In some cases the effect due to the error recovering power prevails on the increment of the base learner error (Figs 8(c), (d)); in other cases, the error recovering power is counter-balanced by the increased average base learner error (Fig. 8(f)); in other cases we can also have similar trends of the ensemble error and the average base learner

error with respect to MHD (Figs 9(b), (c)). We can also observe in a few cases that the average base learner error decreases and the ensemble error increases for some values of the MHD (Figs 9(a), (c)). This seems to be counter-intuitive, but we know that the average base learner error does not take into account the distribution of the error among the dichotomies: we can have very different distributions of the error (and very different ensemble errors) with the same average base learner error. Table 3 shows that there is a relatively large variability of the base learner error.

Hence our results show that there is not a simple relationship between error recovering power and accuracy of the induced dichotomies with respect to the accuracy of the ensemble. It is not sufficient to increase the MHD to improve the performance: we have to take into account the accuracy of the dichotomisers and the dependency among the codeword bit errors.

Following the approach of Ghani [26], with some strong assumptions, we can try to model ECOC ensemble error through the binomial distribution. In fact if we assume that the probability of error of each dichotomiser has about the same value p , and assuming that the outcomes of the n dichotomisers of the ECOC ensemble are independent, the probability of k errors on n trials (dichotomiser tasks) with equal probability of error p is distributed according to a binomial distribution:

$$P(k, n, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

Considering that ECOC codes with $MHD = \min$ can correct $(\min - 1)/n$ errors (Eq. 6), then the overall probability of error of an ECOC ensemble with average base learner error p with $k\min = \frac{\min - 1}{n} + 1$, is given by

$$P_{err} = \sum_{k=k\min}^n \binom{n}{k} p^k (1-p)^{n-k} \quad (11)$$

Table 4 compares the experimental error of the ECOC ensemble and the theoretical error estimated through the binomial distribution (Eq. 11). The results show that only in some cases the experimental and the theoretical error agree: especially for relatively low Hamming distance between codeword the difference is very significant. This is not too surprising, because our assumption about the independence of codeword bit errors and the equal probability of error for all the dichotomisers are too strong. In fact our previous results (Sect. 3.2) showed that a correlation between codeword bit errors does exist, depending on the type of the decomposition, the base learner used and the structure of the data, and also the probability of errors of the different dichotomisers shows a not negligible variance (Table 3). These assumptions can or cannot hold depending on the structure of the data, the type of decomposition and the base learner used.

An important item only indirectly studied in this work is the relationship between the error recovering power and the complexity of the induced dichotomies. Although the average base learner error can be informative about the complexity of the dichotomies, this information is biased by the characteristics of the base learner: the functions implemented by a particular base learner can perform better on certain data distributions, but can achieve worse performances on others [51].

We need more studies to relate the accuracy of the ECOC ensemble with the complexity of the induced decomposition. The relationship effectiveness of ECOC ensemble-complexity of the data is an item common to other ensemble methods [52] and require specific studies and experimental analysis using appropriate measures of complexity, based on geometrical or topological characteristics of data [53].

Conclusions

The effectiveness of ECOC methods depends on many factors ranging from the architecture of the decomposition

Table 3 Average base learner error and its variation. The last two columns refer to the minimum and average diversity of the dichotomies. MHD stands for Minimum Hamming Distance, Stdev for standard deviation, MRD for Minimum Row Distance and ARD for Average Row Distance

Data set	MHD	Average base learner err.	Stdev base learner err.	Min. base learner err.	Max. base learner err.	MRD	ARD
<i>p20</i> 50-bit ECOC	9	0.0912	0.0425	0.0358	0.2675	2	10.075
	11	0.0846	0.0299	0.0358	0.1641	4	9.984
	13	0.0831	0.0260	0.0325	0.1483	4	9.969
	15	0.0816	0.0327	0.0291	0.2150	3	10.030
	17	0.0801	0.0276	0.0325	0.1650	1	10.048
	18	0.0857	0.0255	0.0300	0.1366	3	9.958
	19	0.0939	0.0323	0.0358	0.1908	3	10.018
	20	0.0942	0.0321	0.0358	0.1908	3	9.994
<i>optdigits</i> 64 bit ECOC	15	0.1094	0.1292	0.0042	0.5592	1	4.979
	17	0.1125	0.1296	0.0071	0.5368	1	4.987
	19	0.1336	0.1381	0.0133	0.5386	1	4.959
	21	0.1332	0.1512	0.0071	0.5699	1	4.997
	23	0.1220	0.1281	0.0142	0.5686	1	5.038
	25	0.1473	0.1566	0.0138	0.5919	1	4.991
	27	0.1476	0.1543	0.0090	0.5238	1	4.968
	28	0.1514	0.1557	0.0138	0.5339	1	5.025

Table 4 Experimental ensemble error and theoretical ensemble error predicted through the binomial distribution

Data set	Mi Hamm.dist.	Average base learner err.	Experimental err.	Theoretical err.
<i>p20</i> 50-bit ECOC	9	0.0912	0.0808	0.4846
	11	0.0846	0.0708	0.2462
	13	0.0831	0.0800	0.1187
	15	0.0816	0.0766	0.0483
	17	0.0801	0.0633	0.0167
	18	0.0857	0.0741	0.0249
	19	0.0939	0.0833	0.0165
	20	0.0942	0.0816	0.0168
<i>optdigits</i> 64 bit ECOC	15	0.1094	0.0674	0.4022
	17	0.1125	0.0695	0.2910
	19	0.1336	0.0750	0.3493
	21	0.1332	0.0651	0.2278
	23	0.1220	0.0663	0.0847
	25	0.1473	0.0725	0.1401
	27	0.1476	0.0660	0.0817
	28	0.1514	0.0790	0.0959

unit, to the dependence among codeword bits coding the classes, the decoding function selected for the decision unit, the error recovering power of the ECOC codes, the type and accuracy of the base learners of the ensemble, the complexity of the dichotomies induced by the ECOC decomposition.

The results of our experiments suggest that ensembles of learning machines achieve in general better results than single monolithic learning machines, as dedicated and independent base learners reduce the correlation among the codeword bits, and their learning tasks are reduced to dichotomies in general simpler than polychotomies.

The dependence among codeword bits reduces the error recovering power of ECOC: improving the diversity of the dichotomies and of the dichotomisers can enhance the performance of ECOC learning machines.

Using Minkowski norm in decoding functions of the decision unit seems to be more reliable and robust than using the Hamming distance.

The selection of proper base learners influences the accuracy of the ensemble: a possible way of research could be to experiment with different and specific dichotomic learning machines well-suited for each different dichotomic problem induced by the decomposition.

ECOC codes can recover errors committed by the base learners, but increasing the minimum Hamming distance between codewords does not lead by itself to better performances, because we could have a contemporary increased complexity of the induced dichotomies, or more similar and correlated dichotomies.

Increasing the codeword length, as shown in [9,26] can in general lead to better performances, but for fixed-length codewords our experimentation showed that many factors interact to determine the effectiveness of ECOC methods. We know that the problem of finding an optimal decomposition matrix is NP-complete, and our experimental analysis showed that no straightforward solution exists to select jointly low-correlated dichotomies, code-

words with high minimum Hamming distance and simple induced dichotomies, and base learner well-suited for a given decomposition. The main problem arising from our experimental analysis consist in evaluating how the complexity of the data characterising a given classification problem affects the performance of ECOC methods with respect to the error recovering power of the ECOC decomposition and the accuracy of the base learner of the ensemble. This in turn requires the definition and usage of proper data complexity measures, such as the length of class boundary [54] or the ϵ -neighborhoods space covering [55] to characterise the classification problem complexity on the basis of the geometrical and topological properties of the classes.

Acknowledgements

We thank the anonymous reviewers for their comments and suggestions. This work was supported by INFN, University of Genova, and Madess II CNR.

References

- 1 Dietterich TG (2000) Ensemble methods in machine learning. In: Multiple Classifier Systems. First International Workshop MCS 2000: LNCS 1857, Kittler J, Roli F (eds), Cagliari, Italy, 1–15
- 2 Hashem S (1997) Optimal linear combinations of neural networks. *Neural Computation*, 10: 599–614
- 3 Perrone MP, Cooper LN (1993) When networks disagree: ensemble methods for hybrid neural networks. In: *Artificial Neural Networks for Speech and Vision*, Mammone RJ (ed), 126–142. Chapman & Hall, London
- 4 Freund Y, Schapire RE (1996) Experiments with a new boosting algorithm. In: *Proceedings 13th International Conference on Machine Learning*, 148–156. Morgan Kaufman
- 5 Schapire RE (1999) A brief introduction to boosting. In: *16th International Joint Conference on Artificial Intelligence*, Dean T (ed), 1401–1406. Morgan Kaufman

- 6 Schapire RE, Freund Y, Bartlett P, Lee W (1998) Boosting the margin: A new explanation for the effectiveness of voting methods. *Ann Stat* 26(5): 1651–1686
- 7 Breiman L (1996) Bagging predictors. *Machine Learn* 24(2): 123–140
- 8 Jacobs RA (1995) Methods for combining experts probability assessment. *Neural Computation*, 7: 867–888
- 9 Jordan MI, Jacobs RA (1994) Hierarchical mixture of experts and the EM algorithm. *Neural Computation*, 6: 181–214
- 10 Cherkauer KJ (1996) Human expert-level performance on a scientific image analysis task by a system using combined artificial neural networks. In: Working notes of the AAAI Workshop on integrating Multiple Learned Models, Chan P (ed), 15–21
- 11 Dietterich TG, Bakiri G (1991) Error-correcting output codes: A general method for improving multiclass inductive learning programs. In: Proceedings AAAI-91, 572–577. AAAI Press/MIT Press
- 12 Dietterich TG, Bakiri G (1995) Solving multiclass learning problems via error-correcting output codes. *J Artif Intell Res* (2): 263–286
- 13 James G, Hastie T (1998) The error coding method and pict. *J Computational Graphical Statistics*, 7(3): 377–387
- 14 Kong BE, Dietterich TG (2000) Probability estimation via error correcting output coding. In: IASTED International Conference: Artificial Intelligence and Soft Computing, Banff, Canada
- 15 Kong E, Dietterich TG (1995) Error-correcting output coding correct bias and variance. In: The XII International Conference on Machine Learning, 313–321, San Francisco, CA
- 16 Bose RC, Ray-Chaudhuri DK (1960) On a class of error correcting binary group codes. *Information and Control* (3): 68–79
- 17 Lin S, Costello Jr DJ (1983) Error Control Coding: Fundamentals and Applications. Prentice-Hall, Englewood Cliffs
- 18 Sejnowski TJ, Rosenberg CR (1987) Parallel networks that learn to pronounce english text. *J Artif Intell Res* (1): 145–168
- 19 Schapire RE (1997) Using output codes to boost multiclass learning problems. In: Proceedings Fourteenth International Conference on Machine Learning, 313–321 San Francisco, CA
- 20 James G (1998) Majority vote classifiers: theory and applications. PhD thesis, Department of Statistics, Stanford University, Stanford, CA
- 21 Berger A (1999) Error correcting output coding for text classification. In: IJCAI'99: Workshop on Machine Learning for Information Filtering
- 22 Ghaderi R, Windeatt T (2000) Circular ECOC, a theoretical and experimental analysis. In: Int Conf. Pattern Recognition Barcelona, Spain, 203–206
- 23 Windeatt T, Ghaderi R (2001) Binary labelling and Decision Level Fusion. *Infor Fusion*, 2(2): 103–112
- 24 Allwein EL, Schapire RE, Singer Y (2000) Reducing multiclass to binary: a unifying approach for margin classifiers. *J Machine Learn Res* 1: 113–141
- 25 Aha D, Bankert R (1997) Cloud classification using error-correcting output codes. In: Artificial Intelligence Applications: Natural Science, Agriculture and Environmental Science, vol. 11, 13–28
- 26 Ghani R (2000) Using error correcting output codes for text classification. In: *JCML 2000: Proceedings of the 17th International Conference on Machine Learning*, 303–310, San Francisco, CA
- 27 Bakiri G, Dietterich TG (1999) Achieving high accuracy text-to-speech with machine learning. In *Data mining in speech synthesis*
- 28 Pardo M, Sberveglieri G, Masulli F, Valentini G (2001) Decompositive classification models for electronic noses. *Anal. Chimica Acta*, 446: 223–232
- 29 Kittler J, Ghaderi R, Windeatt T, Matas G (2001) Face Verification using Error Correcting Output Codes. In: *Computer Vision and Pattern Recognition CVPRO1*, Hawaii, USA, 755–760
- 30 Valentini G (2001) Classification of human malignancies by machine learning methods using DNA microarray gene expression data. In: *Fourth International Conference Neural Networks and Expert Systems in Medicine and HealthCare*, Papadourakis GM (ed), 399–408, Milos island, Greece
- 31 Valentini G (2000) Upper bounds on the training error of ECOC-SVM ensembles. Technical Report TR-00-17, DISI, Dipartimento di Informatica e Scienze dell' Informazione, Università di Genova, 2000. <ftp://ftp.disi.unige.it/person/ValentiniG/papers/TR-00-17.ps.gz>
- 32 Masulli F, Valentini G (2000) Effectiveness of error correcting output codes in multiclass learning problems. *Lecture Notes in Computer Science* 1857, 107–116
- 33 Mayoraz E, Moreira M (1997) On the decomposition of polychotomies into dichotomies. In: *The XIV International Conference on Machine Learning*, 219–226, Nashville, TN
- 34 Alpaydin E, Mayoraz E (1999) Learning error-correcting output codes from data. In: *ICANN'99*, Edinburgh, UK, 743–748
- 35 Crammer IC, Singer Y (2000) On the learnability and design of output codes for multiclass problems. In: *Proceedings Thirteenth Annual Conference on Computational Learning Theory*, 35–46
- 36 Masulli F, Valentini G (2000) Comparing decomposition methods for classification. In: *KES'2000, Fourth International Conference on Knowledge-Based Intelligent Engineering Systems & Allied Technologies*, Howlett RJ, Jam LC (eds), 788–791, Piscataway, NJ
- 37 Peterson WW, Weldon Jr EJ (1972) Error correcting codes. MIT Press, Cambridge, MA
- 38 Van Lint J (1971) Coding theory. Springer Verlag, Berlin
- 39 Cortes C, Vapnik V (1995) Support vector networks. *Machine Learn* 20: 273–297
- 40 Alpaydin E, Mayoraz E (1998) Combining linear dichotomisers to construct nonlinear polychotomisers. PAA review Technical Report IDIAP-RR 98-05, IDIAP, Dalle Molle Institute for Perceptual Artificial Intelligence, Martigny
- 41 Masulli F, Valentini G (2000) Parallel Non linear Dichotomisers. In: *JCNN2000, The IEEE-INNS-ENNS International Joint Conference on Neural Networks*, vol. 2, 29–33, Como, Italy
- 42 Nilsson NJ (1965) Learning Machines. Mc Graw Hill, New York
- 43 Anand R, Mehrotra G, Mohan GK, Ranka S (1995) Efficient classification for multiclass problems using modular neural networks. *IEEE Trans Neural Net* 6: 117–124
- 44 Merz CJ, Murphy PM (1998) UCI repository of machine learning databases. www.ics.uci.edu/mllearn/MLRepository.html
- 45 Valentini G, Masulli F (2001) NEUROObjects: an object-oriented library for neural network development. *Neurocomputing* 55(3–4) 623–646
- 46 Everitt BS (1977) The analysis of contingency tables. Chapman & Hall, London
- 47 Dietterich TG (1998) Approximate statistical test for comparing supervised classification learning algorithms. *Neural Computation*, 10:1895–1923
- 48 Guruswami V, Sahai A (1999) Multiclass learning, boosting, and error-correcting codes. In: *Proceedings Twelfth Annual Conference on Computational Learning Theory*, 145–155
- 49 Masulli F, Valentini G (2001) Quantative Evaluation of Dependence among Outputs in ECOC Classifiers Using Mutual Information Based Measures. In: Marko K, Webos P (eds), *Proceedings International Joint Conference on Neural Networks IJCNN'01*, vol. 2, 784–789, Piscataway, NJ
- 50 Masulli F, Valentini G (2001) Mutual information methods for evaluating dependence among outputs in learning machines. Technical Report TR-01-02, DISI, Dipartimento di Informatica e Scienze dell' Informazione, Università di Genova. <ftp://ftp.disi.unige.it/person/ValentiniG/papers/TR-01-02.ps.gz>
- 51 Cohen S, Intrator N (2001) Automatic Model Selection in a Hybrid Perceptron/Radial Network. In: *Multiple Classifier Systems. Second International Workshop, MCS 2001: LNCS 2096*, Cambridge, UK, 349–358
- 52 Ho TK (2001) Data Complexity Analysis for Classifiers Combination. In: *Multiple Classifier Systems. Second International Workshop, MCS 2001: LNCS 2096*, Kittler J, Roli F (eds), Cambridge, UIC, 53–67
- 53 Li M, Vitanyi P (1993) An Introduction to Kolmogorov Complexity and its Applications. Springer-Verlag, Berlin
- 54 Friedman JH, Rafsky LC (1979) Multivariate generalizations of the Wald-Wolfowitz and Smirnov two-sample tests. *Annf Statistics* 7(4): 697–717
- 55 Lebourgeois F, Emptoz H (1996) Pretopological approach for supervised learning. In: *13th Int. Conf. Pattern Recognition*, Wien, Austria, 256–260

Francesco Masulli is an Associate Professor of Computer Science with the University of Pisa (Italy). He authored or co-authored more than 100 papers on Machine Learning, Neural Networks, Fuzzy Systems and Ensemble Methods and co-edited three books and two

special issues of scientific journals on those subjects. He serves as an Associate Editor the international journal *Intelligent Automation and Soft Computing*. He chaired the Conference of the International Graphonomics Society (IGS) in 1997, and the Symposium on Soft Computing SOCO, in 1999. He is a member of the IEEE-Neural Network Council (Italian R.I.G.), and a Board Member of the Italian Neural Network Society (SIREN) and of the SIG Italy of the International Neural Network Society (INNS).

Giorgio Valentini is research assistant with the Department of Computer Science (DSI), University of Milan (Italy). He received the Ph.D. in Computer Science from the University of Genova. He is member of the International Society of Computational Biology (ISCB), INNS and SIREN. His research interests include ensembles of learning machines and bioinformatics.