

Error-Correcting Classification via ML-Techniques and Regularization

Wolfgang Utschick and Josef A. Nossek

Technische Universität München
 Lehrstuhl für Netzwerktheorie und Schaltungstechnik
 Arcisstr. 21, D-80290 München
 wout@nws.e-technik.tu-muenchen.de

Abstract—Instead of *a priori* approaches to Error-Correcting Output Coding (ECOC) an adaptive solution to the definition of codewords in multiclass learning problems is presented. The generation of codewords and the training of parallel subsystems of the overall classification system is matter of an iterative optimization method which takes advantage of a Maximum Likelihood (ML) technique in combination with a regularization principle for penalizing low distances between codewords. Even for a very restricted complexity of outputs the optimized codes outperform the trivial one-per-class output coding. In a practical application, the generalization capability of the inferred decision rule was considerably improved.

I. Introduction

Polychotomous classification is a general task in pattern recognition. Designing a classifier is considered to be the approximation of an unknown decision rule which will be applied to a sequence of patterns. The input pattern must be uniquely assigned to one element of the unordered set of class indices $\mathcal{K} = \{1, 2, \dots, K\}$ on the basis of n observed attributes $(x_1, \dots, x_N) = \mathbf{x} \in \mathcal{R}^N$ from the original pattern. The construction of the decision rule is inferred from a set of pre-classified training examples $\mathcal{S} = ((\mathbf{x}_1, k_1), (\mathbf{x}_2, k_2), \dots, (\mathbf{x}_M, k_M)) \subset \mathcal{X} \times \mathcal{K}$, for which the true class indices $k_m \in \mathcal{K}$ are given (see supervised learning in [1, 2, 3, 4]).

The approximation of the decision rule is generally identified with the construction of a mapping from the feature space \mathcal{R}^N into the decision space \mathcal{D} [3]. The decision process of each input vector is composed in two steps by

$$d: \mathcal{R}^n \rightarrow \mathcal{D} \rightarrow \mathcal{K}, \mathbf{x} \mapsto d \mapsto k. \quad (1)$$

A highly important feature between the layers of the classifier is the representation of classes in the decision space. Practitioners often prefer the trivial 1-out-of- K coding where each output of the classifier corresponds to one of all classes. This method directly corresponds

to decision rules based on posterior probabilities of classes [5]. Provided the representation of classes is based on binary reference vectors $\mathbf{t} \in \{-1, +1\}^J$ embedded in a real-valued decision space \mathcal{D} , the estimation of reference vectors (codewords) for classes [6] is equal to the optimal decomposition of polychotomies into dichotomies [7] (e.g. class k is associated with one of the two super classes of the j th dichotomy if $t_{kj} = +1$, and v.v.). Reference vectors can be interpreted as corners of a high-dimensional hyper cube.

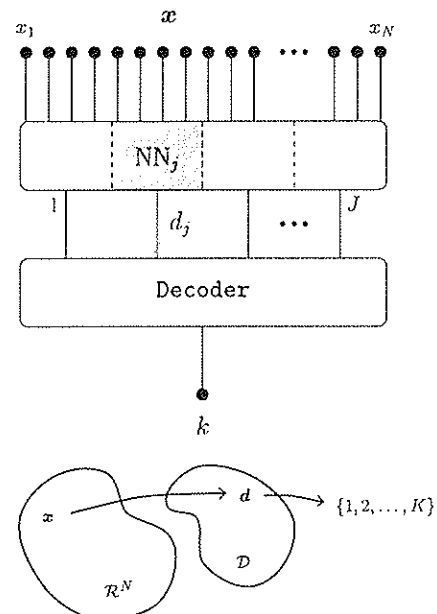


Figure 1: The Classification System.

It has been shown that ECOC proposed by T. Dietterich and G. Bakiri [6, 8] outperforms the one-per-class approach. ECOC is a kind of “neither homogeneous nor non-homogeneous voting” [8], which reduces both the bias error and the variance error if the bias errors between the individual two-group classifiers

are uncorrelated. On the other hand, error-correcting coding is supposed to be robust in relation to small sample sizes and in some cases does not depend on the particular assignment of reference vectors to classes [8]. Whereas T. Dietterich's method of ECOC is *a priori*, the decomposition of polychotomies proposed by E. Mayoraz and M. Moreira [7] depends on the underlying training data.

In this work we propose a method that estimates the reference vectors of classes using a technique that combines a maximum likelihood (ML) method for density estimation together with the principle of regularization. The presented classification system (Fig. 1) consists of J parallel classifiers (neural networks) and a following decoder that realizes the function from the decision space onto the set of indices. Each parallel network realizes its own two-class problem with respect to the decomposition of the complete classification problem. The $\mathbf{d}(\mathbf{x}, \boldsymbol{\theta})$ expresses the mapping of the first layer where $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_J)$ summarizes the parameters of the parallel components, and $\mathbf{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_K)^T$ represents the set of desired codewords for classes in form of a matrix. The decoder performs a minimum distance based decision in form of

$$k = \arg \min_{l \in K} \|\mathbf{d} - \mathbf{t}_l\|_2. \quad (2)$$

Instead of two separate optimization tasks – searching for a set of reference vectors (decomposition into dichotomies) and training of the parallel components – the training of the network layer is split into a number of training epochs. During each training epoch the parameters of the network components and the reference vectors of classes are iteratively estimated by means of expectation maximization (EM) [9, 10].

In contrast to EM in density estimation problems, the transformed input vectors in form of $\mathbf{d}(\mathbf{x}, \boldsymbol{\theta})$, or in other words the realizations of the random variable $\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_M$ – the given sample which plays the fundamental part in the ML algorithm – are not constant. Therefore the algorithm rather aims at learning the parameters $\boldsymbol{\theta}$ of all parallel classifiers such that the vectors $\mathbf{d}_m(\boldsymbol{\theta}) = \mathbf{d}(\mathbf{x}_m, \boldsymbol{\theta})$, $\mathbf{x}_m \in \mathcal{X}$, optimally correspond to mixtures of disjoint distributions $p(\mathbf{d}(\mathbf{x}, \boldsymbol{\theta}) | \mathbf{t})$ according to their respective class memberships. The estimation of this mixture distributions is in turn part of the optimization. It is assumed that for a training epoch the representation of each class is distributed over a small subset of vectors $\mathbf{t} \in \mathcal{T}_k$. The subset contains neighboring vectors of reference vector \mathbf{t}_k which is defined to be the “gravity center” of all $\mathbf{t} \in \mathcal{T}_k$ during this epoch. In order to prevent the distributions of classes in decision space from overlapping, the subsets of vectors $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_K$ must be disjoint. Under the assumption of normal distributions around vectors $\mathbf{t} \in \mathcal{T}_k$, the center of distributions are impli-

city given. The covariance of the normal distribution is defined to be diagonal and only depends on a single variance parameter which is reasonably set to > 1 (cf. distance between binary reference vectors).

The characteristic of the EM algorithm is that the optimization of both objectives, namely the optimum mixture of distributions according to $\mathbf{t} \in \mathcal{T}_k$ and the embedding of training samples into desired distributions by training of the first layer, is alternately performed. The algorithm is only indirectly error-correcting! Fig. 2 presents the idea of the algorithm.

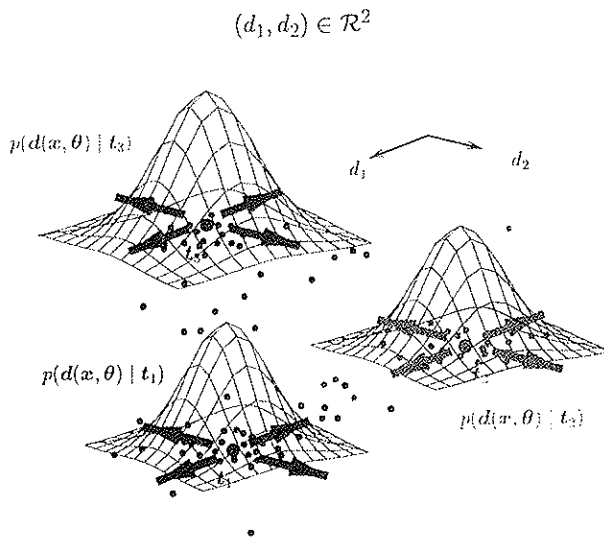


Figure 2: “Gravity Centers” in Decision Space.

At the end of each training epoch, the set of reference vectors \mathbf{T} is updated. The update of reference vectors allows to implement the principle of regularization [11]. Using a regularization term, which is based on distances between reference vectors in decision space, already have shown promising results [12].

II. Expectation Maximization

Given the estimated parameters $\boldsymbol{\theta}^*$ and \mathbf{T}^* ($P^*[\mathbf{t}]$) from the last training epoch, the expectation term (cf. Appendix) of the EM is equal to

$$E = \int_{\mathcal{T}_{k_1}^*} \dots \int_{\mathcal{T}_{k_M}^*} \sum_{m=1}^M \log[p(\mathbf{d}_m(\boldsymbol{\theta}), \mathbf{t}_m)] \prod_{\mu=1}^M p(\mathbf{t}_\mu | \mathbf{d}_\mu(\boldsymbol{\theta}^*)) dt_1 \dots dt_M. \quad (3)$$

The reverse order of intergration and summation in Eqn. 3 together with the condition of

$\int_{\mathcal{T}_{k_\mu}^*} p(\mathbf{t}_\mu | \mathbf{d}_\mu(\theta^*)) d\mathbf{t}_\mu = 1$ results in

$$E = \sum_{m=1}^M \int_{\mathcal{T}_{k_m}^*} \log[p(\mathbf{d}_m(\theta), \mathbf{t})] p(\mathbf{t} | \mathbf{d}_m(\theta^*)) d\mathbf{t}. \quad (4)$$

Considering the prior knowledge of class membership of vectors \mathbf{d}_m , which reduces the variety of $\mathcal{T}_{k_m}^*$ to K disjoint \mathcal{T}_k^* , and $\log[p(\mathbf{d}_m(\theta), \mathbf{t}_m)] = \log[p(\mathbf{d}_m(\theta) | \mathbf{t}_m)] + \log[p(\mathbf{t}_m)]$, simplifies the E-step to

$$E = \sum_{k=1}^K \int_{\mathcal{T}_k^*} \sum_{m_k=1}^{M_k} \log[p(\mathbf{d}_{km_k}(\theta) | \mathbf{t})] p(\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)) d\mathbf{t} + \sum_{k=1}^K \int_{\mathcal{T}_k^*} \sum_{m_k=1}^{M_k} \log[p(\mathbf{t})] p(\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)) d\mathbf{t} \quad (5)$$

subject to

$$\int_{\mathcal{T}_k^*} p(\mathbf{t}) d\mathbf{t} = P[k], \quad (6)$$

where $P[k]$ denotes the prior probability of class index k , and $M_1 + M_2 + \dots + M_K = M$.

Finally, the next estimation of parameters is obtained by maximizing both expectation terms in Eqn. 5. Using binary vectors $\mathbf{t}_k \in \mathcal{T}_k^* \subset \{-1, +1\}^J$, \int is replaced by \sum , $p(\mathbf{t}) \rightarrow P[\mathbf{t}]$, and each domain \mathcal{T}_k^* contains a finite number of possible reference vectors for class k .

For the maximization of the expectation term (5), the gradient according to θ is

$$\frac{\partial E}{\partial \theta} = \sum_{k=1}^K \sum_{\mathbf{t}} \sum_{m_k=1}^{M_k} \frac{P[\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)]}{p(\mathbf{d}_{km_k}(\theta) | \mathbf{t})} \frac{\partial p(\mathbf{d} | \mathbf{t})}{\partial \mathbf{d}} \frac{\partial \mathbf{d}}{\partial \theta}. \quad (7)$$

Analogously,

$$\frac{\partial E}{\partial P[\mathbf{t}]} = \sum_{m_k=1}^{M_k} \frac{P[\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)]}{P[\mathbf{t}]} + \text{const.} \quad (8)$$

for all $\mathbf{t} \in \bigcup_{k \in \mathcal{K}} \mathcal{T}_k^*$, where *const.* represents the Lagrangian multiplier from the additional linear constraint $P[k=1] + P[k=2] + \dots + P[k=K] = 1$ together with Eqn. 6. Setting Eqn. 8 to zero, the summation over all $\mathbf{t} \in \mathcal{T}_k^*$ yields *const.* = $-M_k$, and for all $k \in \mathcal{K}$ and $\mathbf{t} \in \mathcal{T}_k^*$ the estimation of $P[\mathbf{t}]$ is obtained by

$$P[\mathbf{t}] = \frac{1}{M_k} \sum_{m_k=1}^{M_k} P[\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)], \quad (9)$$

where $P[\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)]$ is inferred from the estimated parameters of the last training epoch:

$$P[\mathbf{t} | \mathbf{d}_{km_k}(\theta^*)] = \frac{p(\mathbf{d}_{km_k}(\theta^*) | \mathbf{t}) P^*[\mathbf{t}]}{\sum_{\tau} p(\mathbf{d}_{km_k}(\theta^*) | \tau) P^*[\tau]}. \quad (10)$$

III. Regularization

In order to prevent the distributions of classes in decision space from overlapping, and for the benefit of ECOC, increasing the minimum distance between codewords of classes is highly desirable. In the following, we apply a regularization principle during the learning process for generation of reference vectors. The regularization term is defined as the logarithm of the product of all distances between reference vectors, or

$$R = \sum_{k=1}^{K-1} \sum_{l=k+1}^K \log \|\mathbf{t}_k - \mathbf{t}_l\|_2, \quad (11)$$

and favours the generation of higher Hamming distances between reference vectors. Hence, the cost function of the overall optimization is defined as

$$F = E(\mathcal{X}, \mathcal{K}, \theta, \mathbf{T}^*) + \lambda \cdot R(\mathbf{T}^*), \quad (12)$$

where λ is called "regularization parameter".

After each training epoch the next potential reference vector of class k is calculated by

$$\mathbf{t}_k = \mathbf{t}_k^* \oplus (\nabla_k E + \lambda \cdot \nabla_k E). \quad (13)$$

The ∇_k is called pseudo-gradient operator because in the theory of the differential calculus there is no definition of gradients in $\{-1, +1\}^J$. The definition in Eqn. 13 roughly follows the rules of numerical differentiation. The feasible candidates for \mathbf{t}_k are elements from the $\mathcal{N}_1(\mathbf{t}_k^*)$ neighborhood ($\mathbf{t} \in \mathcal{T}_k^*$ which maximally differ from \mathbf{t}_k^* in one component) of the current reference vector. The operator \oplus is defined as $\mathbf{t}_k \oplus \mathbf{v} := \arg \min_{\mathbf{t} \in \mathcal{T}_k^*} \|\mathbf{t}_k + \mathbf{v} - \mathbf{t}\|_2$ for all $\mathbf{v} \in \mathcal{R}^J$. The restriction to $\mathcal{N}_1(\mathbf{t}_k)$ neighborhoods of possible candidates is important. Otherwise, large modifications of the reference vectors would change the composition of the corresponding dichotomies such dramatically that the after-training meaning of the learned parameters from the previous training epoch would be destroyed.

The components of pseudo-gradient $\nabla_k E$ would be consequently defined by

$$\left. \frac{\partial E}{\partial t_{kj}} \right|_{\mathbf{t}_k^*} = \frac{E(\mathbf{t}_k^* + \epsilon \cdot \mathbf{e}_j) - E(\mathbf{t}_k^*)}{|\epsilon|}, \quad (14)$$

where all e_{ji} of \mathbf{e}_j are equal to δ_{ji} , and $\epsilon = -2t_{kj}^*$ because of the binary properties of codewords. Hence, the calculation of $\nabla_k E$ would require multiple repeated runs of the last training epoch based on slightly modified reference vectors.

A definition of $\nabla_k E$ with lower complexity can be obtained by utilizing the probabilities of vectors $\mathbf{t} \in \mathcal{T}_k^*$ from the previous training epoch (see Eqn. 9). The proposed definition is equal to

$$\left. \frac{\partial E}{\partial t_{kj}} \right|_{\mathbf{t}_k^*} = \frac{P[\mathbf{t}_k^* + \epsilon \cdot \mathbf{e}_j] - P[\mathbf{t}_k^*]}{|\epsilon|}. \quad (15)$$

The pseudo-gradient of the regularization term is analogously derived from the rules of numerical differentiation:

$$\begin{aligned} \frac{\partial}{\partial t_{kj}} R \Big|_{\mathbf{t}_k^*} = & \quad (16) \\ & \sum_{l=1}^{k-1} \frac{1}{\|\mathbf{t}_k - \mathbf{t}_l\|_2} \cdot \frac{\partial}{\partial t_{kj}} \|\mathbf{t}_k - \mathbf{t}_l\|_2 \Big|_{\mathbf{t}_k^*} \\ & + \sum_{l=k+1}^K \frac{1}{\|\mathbf{t}_k - \mathbf{t}_l^*\|_2} \cdot \frac{\partial}{\partial t_{kj}} \|\mathbf{t}_k - \mathbf{t}_l^*\|_2 \Big|_{\mathbf{t}_k^*}. \end{aligned}$$

The differential ($k \neq l$) is again defined as

$$\begin{aligned} \frac{\partial}{\partial t_{kj}} \|\mathbf{t}_k - \mathbf{t}_l^*\|_2 \Big|_{\mathbf{t}_k^*} = & \quad (17) \\ \frac{\|\mathbf{t}_k^* + \epsilon \cdot \mathbf{e}_j - \mathbf{t}_l^*\|_2 - \|\mathbf{t}_k^* - \mathbf{t}_l^*\|_2}{|\epsilon|}. \end{aligned}$$

In order to guarantee that $F(\mathcal{X}, \mathcal{K}, \boldsymbol{\theta}, \mathbf{T}) \geq F(\mathcal{X}, \mathcal{K}, \boldsymbol{\theta}^*, \mathbf{T}^*)$, the iterative procedure (13) of the discrete optimization must be performed semi-implicitly and the next set of possible candidates of reference vectors \mathbf{T} must be chosen properly. This is because of the non-local properties of the pseudo-gradients. The calculation of $\nabla_k R$ for the k -th reference vector depends on the already-adapted reference vectors $\mathbf{t}_1 \dots \mathbf{t}_{k-1}$, and the still-unadapted vectors $\mathbf{t}_k^* \dots \mathbf{t}_K^*$. The order of the sequence of updates influences the optimization result of the iteration step. A random choice of the order avoids systematical errors.

If the pseudo-gradient points into the direction of a possible candidate for \mathbf{t}_k , it simply depends on a threshold whether the corresponding bit of the original reference vector changes from $+1 \rightarrow -1$ or $-1 \rightarrow +1$, or does not change. The threshold is practically determined by a stepsize.

There might be an interesting interpretation of the proposed cost function approach: structural risk minimization (SRM) and regularization. SRM [13] is a result of Statistical Learning Theory and emphasizes the minimization of the classification error on training samples as well as the generalization on test samples by optimal model selection. The impact of the regularizer R complicates the generation of trivial two-class problems (dichotomies) of each parallel classifier (the simplest dichotomy would consist of purely one class separated from all other classes). Hence, regularization increases the expressivity requirements which are imposed on the selected model by the composed classification task. This may be interpreted as a type of SRM by configuration of the training data and not by model selection.

IV. Example

The following example is based on the data sets from a realistic practical application – the recognition of handwritten digits ($K = 10$) taken from a mixture of the NIST databases [14]. The choice of the classification task is especially relevant because of its importance on the generation of optimal dichotomies according to their data structure. The training was not applied to the pure bitmap of scanned images but to the transformed representation of the data. Each image of a character was transformed into a high-dimensional vector ($N = 194$) by a feature extraction method [15]. The feature extraction is performed by matching the original input images with a variety of artificial patterns, and each of this patterns is generated as a grid of equidistant lines or concentric circles. The matching scores are equal to the input vector \mathbf{x} .

In a very first trial, the expectation term (5) was substituted by means of the empirical error of the decision rule (see [12]). Then

$$E \approx \sum_{k=1}^K P[\mathbf{d}(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{t}_k \mid \mathbf{x} \in \mathcal{X}_k] \cdot P[k] \quad (18)$$

is inspired from the supposed similarity between maximizing E and minimizing the empirical error of the decision rule under the assumption of error-correcting training of the classifier, and $\mathcal{T}_k = \mathcal{N}_e(\mathbf{t}_k)$ where $e = \lfloor \frac{h_{min}}{2} \rfloor$, and h_{min} expresses the minimal Hamming distance between any codewords of \mathcal{T} .

The following results are based on the minimization of the empirical error instead of applying Eqn. 7. The employed learning algorithm is an extended version of the classical error-correcting Madaline algorithm [16] (see [12]) which is separately applied to the single dichotomies of each parallel network and exclusively depends on the binary reference values $t_{kj} \in \{-1, +1\}$ of components. The one per class coding (1-out-of-10) served as a startcode for the iterative optimization process. The number of parallel components is equal to the number of classes ($J = 10$). Each parallel network consists of 5 neurons and the number of free parameters is equal to $5 \times 194 \approx 1000$. The sample size of \mathcal{S} and the size of an i.i.d. set of test examples is 10×2000 respectively.

The influence of the regularization parameter λ is presented in Figure 3. It shows averaged results (incl. standard deviation) from 10 repeated training processes. The error rate of the constant 1-out-of-10 coding serves as a reference. The polynomial least squares approximation of the error rates versus λ is given by $\lambda \mapsto 2.5707 - 0.9489 \cdot \lambda + 0.6436 \cdot \lambda^2$. For the given configurations the optimal value of parameter was found to be $\lambda_{opt} \approx 0.8$. For larger values of λ , the generation of reference vectors would have resulted in ma-

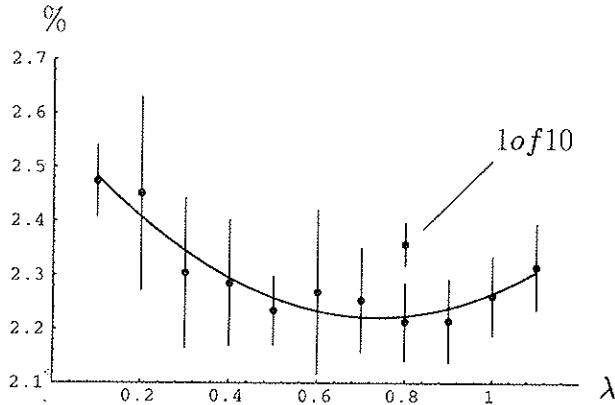


Figure 3: Error Rates of Test Samples after Training vs. Regularization Parameter.

ximal Hamming distances, however almost independently from the underlying database, and the generalization capabilities would have become worse.

Eqn. 19 presents the codematrix after 18 training epochs of 2000 online parameter adaptations of θ respectively:

$$\mathbf{T}_{opt} = \begin{pmatrix} +1 & -1 & +1 & -1 & -1 & -1 & -1 & +1 & +1 & -1 \\ -1 & +1 & -1 & -1 & +1 & -1 & -1 & -1 & +1 & +1 \\ -1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 & -1 & -1 \\ +1 & -1 & -1 & +1 & -1 & -1 & +1 & -1 & -1 & +1 \\ -1 & -1 & -1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 \\ -1 & +1 & -1 & +1 & -1 & +1 & -1 & +1 & -1 & -1 \\ +1 & +1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 & -1 \\ +1 & +1 & -1 & -1 & +1 & -1 & -1 & +1 & -1 & -1 \\ +1 & -1 & -1 & +1 & -1 & +1 & -1 & -1 & +1 & -1 \\ -1 & -1 & +1 & -1 & -1 & +1 & -1 & +1 & -1 & +1 \end{pmatrix}. \quad (19)$$

Although the probability of error-free decisions – correct decisions at any component – became even worse, the generalization error of the decision rule – exploiting the error-correcting properties of \mathbf{T}_{opt} – was reduced about 12% from 2.35% (1-out-of-10) to 2.07%. The training results of the overall classification rule and its components have always met 0.00% classification error.

$$\{h_{kl}\} = \begin{pmatrix} 6 & 0 & 5 & 6 & 7 & 6 & 4 & 4 & 4 & 4 \\ 6 & 5 & 5 & 6 & 5 & 6 & 6 & 4 & 6 & 6 \\ 5 & 5 & 5 & 5 & 4 & 5 & 5 & 5 & 5 & 5 \\ 6 & 0 & 5 & 5 & 5 & 6 & 4 & 6 & 4 & 6 \\ 7 & 5 & 4 & 5 & 5 & 5 & 5 & 5 & 5 & 5 \\ 6 & 6 & 5 & 6 & 5 & 6 & 4 & 4 & 4 & 4 \\ 4 & 6 & 5 & 4 & 5 & 6 & 4 & 6 & 6 & 6 \\ 4 & 4 & 5 & 6 & 5 & 4 & 4 & 4 & 6 & 6 \\ 4 & 6 & 5 & 4 & 5 & 4 & 6 & 6 & 6 & 6 \\ 4 & 6 & 5 & 6 & 5 & 4 & 6 & 6 & 6 & 6 \end{pmatrix}. \quad (20)$$

The bit-distances between the reference vectors of (19) are displayed in Eqn. 20. Whereas the minimal Hamming distance $h_{min} = 4$, the averaged Hamming distance between reference vectors $\frac{1}{(K-1) \cdot K/2} \sum_{k=1}^{K-1} \sum_{l=k+1}^K h_{kl}$ is even larger than 5. Although the advanced ECOC principally benefits the classification system, the increased correlations bet-

ween errors of different parallel components for non-trivial coding (\neq one per class) remain a problem. Eqn. 23 shows the error correlations between NN_j , $j = 1, 2, 3$ and 4, according to 1-out-of-10 output coding. Analogously, Eqn. 24 holds for the optimized codematrix in Eqn. 19. The correlation matrices are calculated by $\mathbf{R} = \frac{1}{100} \sum_{m=1}^M \frac{1}{M} \mathbf{r}(\mathbf{x}_m) \mathbf{r}(\mathbf{x}_m)^T$ where

$$r_j(\mathbf{x}) := \begin{cases} 1 & , \text{ if } d_j(\mathbf{x}_m) \neq t_{kmj} \\ 0 & , \text{ otherwise} \end{cases}. \quad (21)$$

$$\mathbf{R}_{1of10} = \begin{pmatrix} 0.310 & 0.010 & 0.020 & 0.010 & \dots \\ 0.010 & 0.330 & 0.030 & 0.015 & \dots \\ 0.020 & 0.030 & 0.695 & 0.100 & \dots \\ 0.010 & 0.015 & 0.050 & 0.880 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \quad (22)$$

$$\mathbf{R}_{opt} = \begin{pmatrix} 2.010 & 0.510 & 0.670 & 0.435 & \dots \\ 0.510 & 2.085 & 0.440 & 0.335 & \dots \\ 0.670 & 0.440 & 3.420 & 0.485 & \dots \\ 0.435 & 0.335 & 0.485 & 1.825 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \quad (23)$$

The definition of error-correcting output codes therefore can not be a matter of *a priori* choices. The presented example shows the superiority of optimized reference vectors even for a considerably higher correlation between errors of parallel components and a rather low complexity of outputs. Further experiments confirmed the earlier results [8] that non-trivial coding is more advantageous in cases where overfitting is likely [12].

V. Conclusion

Theoretical results in ECOC as well as practical applications have generally considered the case where the number of outputs is much larger than the number of classes [6, 8, 17]. Results about using error-correcting codes with small code length for classification problems are not known to the author. On the other hand, the use of $J \gg K$ is not recommendable because of the growing size of the classification system and the increased computational costs for the decision process.

We have shown that starting with a trivial 1-out-of- K code, more efficient codes can be produced by iterative optimization. Even for the case of a restricted number of parallel components $J = K$, the one-per-class output coding is not optimal. A more efficient code has been generated by an iterative optimization method based on a maximum likelihood technique combined with a regularization principle for ECOC. The generation of reference vectors is not based on *a priori* appreciated codes but rather depends on an iterative method which considers both the classification

system and the structure of data from the underlying classification task. In the example, the generalization of the classifier for handwritten characters was improved about more than 10%.

Appendix

The EM algorithm is a general method for solving ML estimation problems given incomplete data. Let

$$L(\theta) = \sum_{m=1}^M \log p(\mathbf{x}_m; \theta)$$

denote the log-likelihood function of the training set S where $p(\mathbf{x}; \theta)$, which depends on a set of parameters θ , denotes the density function of the probability of input patterns \mathbf{x} .

After simple manipulations and taking the conditional expectation $[\dots; \theta^*]$ (at parameters θ^*) with respect to suitably introduced data $(\mathbf{y}_m | \mathbf{x}_m)$, one obtains $L(\theta) = E(\theta, \theta^*) + H(\theta, \theta^*)$ where

$$E(\theta, \theta^*) = \left[\sum_{m=1}^M \log p(\mathbf{x}_m, \mathbf{y}_m; \theta) ; \theta^* \right]$$

$$H(\theta, \theta^*) = \left[\sum_{m=1}^M \log p(\mathbf{y}_m | \mathbf{x}_m; \theta) ; \theta^* \right]$$

It can be shown that it suffices to maximize $E(\theta, \theta^*)$ if one aims to maximize the log-likelihood function $L(\theta)$.

The idea of the EM method is to produce a simpler ML estimation problem by any appropriate choice of "complete data". The most important prerequisite for the practicability of the EM algorithm is that the original likelihood function of the "incomplete data" \mathbf{x}_m can be obtained by the likelihood function of the unknown "complete data" $(\mathbf{y}_m | \mathbf{x}_m)$ in form of "marginalization" which involves integrating out the so-called "complete data" [9, 10].

It follows the outline of the EM algorithm: the algorithm starts with an arbitrary initial guess and denotes by θ^* the current estimate of the parameters θ after a number of training epochs. The next iteration step consists of two steps, namely the E-step (expectation):

$$E(\theta, \theta^*) = \int_{\mathcal{Y}_1} \dots \int_{\mathcal{Y}_M} \sum_{m=1}^M \log p(\mathbf{x}_m, \mathbf{y}_m; \theta) \cdot \prod_{m=1}^M p(\mathbf{y}_m | \mathbf{x}_m; \theta^*) d\mathbf{y}_1 \dots d\mathbf{y}_M,$$

and the M-Step (maximization): $\max_{\theta} E(\theta, \theta^*)$.

If $Q(\theta, \theta^*)$ is continuous in both θ and θ^* , it has been proven that the algorithm converges to a local maximum of the log-likelihood function $L(\theta)$ [9, 10].

References

- [1] C. M. Bishop. *Neural Networks for Pattern Recognition*. Clarendon Press, 1995.
- [2] B.D. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.
- [3] J. Schürmann. *Pattern Classification, A Unified View Of Statistical And Neural Approaches*. Wiley, 1996.
- [4] H. Bunke and P.S.P. Wang, editors. *Handbook on Optical Character Recognition and Document Analysis*. World Scientific Publishing Company, 1996.
- [5] S. Geman, E. Bienenstock, and R. Doursat. Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1-58, 1992.
- [6] T.G. Dietterich and G. Bakiri. Solving multiclass learning problems via error-correcting output codes. *Journal of Artificial Intelligence Research*, 2:263-286, 1995.
- [7] E. Mayoraz and M. Moreira. On the decomposition of polychotomies into dichotomies. In *Proceedings of the 14th International Conference on Machine Learning*, pages 219-226, 1996.
- [8] E.B. Kong and T.G. Dietterich. Error-Correcting Output Coding Corrects Bias and Variance. In *Proceedings of the 12th International Conference on Machine Learning*, pages 313-321. Morgan Kaufmann, 1995.
- [9] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, B*, 39(1), 1977.
- [10] R.A. Redner and H.F. Walker. Mixture densities, maximum likelihood and the EM algorithm. *Siam Review*, 26(2):195-239, 1984.
- [11] T. Poggio and F. Girosi. Networks for approximation and learning. *Proceedings of the IEEE*, 78(9):1481-1497, 1990.
- [12] W. Utschick. *Error-Correcting Classification Based on Neural Networks*. PhD thesis, Technische Universität München, 1998. to appear.
- [13] V.N. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- [14] The state-of-the-art in OCR is subject of NIST conference. *Intelligence*, 9(6), 1992.
- [15] W. Utschick, P. Nachbar, C. Knobloch, A. Schuler, and J.A. Nossek. The evaluation of feature extraction criteria applied to neural network classifiers. In *Proceedings of the 3th International Conference on Document Analysis and Recognition*, pages 315-318. IEEE Computer Society Press, 1995.
- [16] B. Widrow and M.A. Lehr. 30 years of adaptive neural networks: perceptron, madaline, and backpropagation. *Proceedings of the IEEE*, 78(9):1415-1441, 1990.
- [17] J. Gareth and T. Hastie. Error coding and PaCT's, 1997. Winning paper in the ASA student paper competition for the Statistical Computing Section, available at <http://playfair.Stanford.EDU/~trevor/ASA/winners.html>.