

Contents lists available at [SciVerse ScienceDirect](#)

Pattern Recognition Letters

journal homepage: www.elsevier.com/locate/patrec

A GA-based feature selection approach with an application to handwritten character recognition

C. De Stefano, F. Fontanella*, C. Marrocco, A. Scotto di Freca

Dipartimento di Ingegneria Elettrica e dell'Informazione (DIEI), Università di Cassino e del Lazio Meridionale, Italy

ARTICLE INFO

Article history:
Available online xxxx

Keywords:
Feature selection
Genetic algorithms
Handwriting recognition

ABSTRACT

In the framework of handwriting recognition, we present a novel GA-based feature selection algorithm in which feature subsets are evaluated by means of a specifically devised separability index. This index measures statistical properties of the feature subset and does not depend on any specific classification scheme. The proposed index represents an extension of the Fisher Linear Discriminant method and uses covariance matrices for estimating how class probability distributions are spread out in the considered N -dimensional feature space. A key property of our approach is that it does not require any a priori knowledge about the number of features to be used in the feature subset. Experiments have been performed by using three standard databases of handwritten digits and a standard database of handwritten letters, while the solutions found have been tested with different classification methods. The results have been compared with those obtained by using the whole feature set and with those obtained by using standard feature selection algorithms. The comparison outcomes confirmed the effectiveness of our approach.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

It is generally agreed that one of the main factors influencing performance in handwriting recognition is the selection of an appropriate set of features for representing input samples (Ahlgren et al., 1971; Shi et al., 1998; Chung et al., 1997; Kim et al., 2000; Oliveira et al., 2003; Nunes et al., 2004; Cordella et al., 2008). This has led to the development of a large variety of feature sets, which are becoming increasingly larger in terms of number of attributes. The aim is to address the problem of diversity in style, size, and shape, which can be found in handwriting produced by different writers (Kim and Govindaraju, 1997). The effect is that the efficiency of learning algorithms may degrade, especially in presence of irrelevant or redundant features.

To overcome this problem and maximize classification performance, many techniques have been proposed for reducing the dimensionality of the feature space in which data have to be processed. These techniques, generally denoted as feature reduction (Fodor, 2002), may be divided in two main categories, called feature extraction and feature selection. Feature extraction-based methodologies transform the original feature space into a smaller one. The transformation can be any linear or nonlinear combination of the original features (Fukunaga, 1990). Feature

selection-based approaches, instead, produce as output a feature subset from the original one, without any kind of transformation (Guyon and Elisseeff, 2003). Such subset is supposed to include the best features according to a certain criterion. The role of such criterion consists in identifying the subset providing the most discriminative power.

The choice of a good feature subset is a crucial step in any classification process for several reasons:

- The features used to describe the patterns determine the search space to be explored during the learning phase. Then, irrelevant and noisy features make the search space larger, increasing both the time and the complexity of the learning process.
- If the considered feature subset does not include all the information needed to discriminate patterns belonging to different classes, the achievable classification performances may be unsatisfactory, regardless the effectiveness of the learning algorithm employed.
- Irrelevant and noisy features improperly chosen may make the learning process ineffective.
- The computational cost of the classification process depends on the number of features used to describe the patterns. Then, reducing such number results in a significant reduction of this cost.

When the cardinality N of the whole feature set Y is high, the problem of finding the optimal feature subset becomes computationally intractable because of the resulting exponential growth

* Corresponding author. Tel.: +39 0776 2993382.

E-mail addresses: destefano@unicas.it (C. De Stefano), fontanella@unicas.it (F. Fontanella), cristina.marrocco@unicas.it (C. Marrocco), a.scotto@unicas.it (A. Scotto di Freca).

of the search space, made of all the 2^N possible subsets of Y . Many heuristic algorithms have been proposed in the literature for finding near-optimal solutions: Greedy selection (Kwak and Choi, 2002), branch and bound (B&B) (Somol et al., 2004), floating search (Somol et al., 1994). These algorithms use greedy stepwise strategies that incrementally generate feature subsets by adding the feature that produces the highest increment of the evaluation function. Since these algorithms do not take into account complex interactions among several features, in most of the cases they lead to sub-optimal solutions. An alternative way to cope with the search problem is that of using genetic algorithms (GAs), which have demonstrated to be an effective search tools for finding near-optimal solutions in complex and non-linear search spaces (Goldberg, 1989). For this reason, GA-based search strategies have been widely used to solve feature selection problems (Kudo and Sklansky, 2000; Oh et al., 2004; Cordella et al., 2010; De Stefano et al., 2007; Siedlecki and Sklansky, 1989; Yang and Honavar, 1998). In (Kudo and Sklansky, 2000; Oh et al., 2004), a Nearest Neighbor (NN) classifier has been used for evaluating feature subsets, while in (Yang and Honavar, 1998) this goal is achieved by using a Neural Network and by combining the classification results with some costs associated to the features. In particular, in (Oh et al., 2004), a hybrid mechanism is proposed for finding better solutions in the neighborhood of each solution found by the GA. Moreover, comparative studies have demonstrated the superiority of GAs in feature selection problems involving large numbers of features (Kudo and Sklansky, 2000). In all the mentioned approaches, however, the cardinality of the subset to be found must be a priori fixed. Finally, in (Chouaib et al., 2008) a GA based method is presented, which uses a combination of Adaboost classifiers for evaluating the fitness of each individual in the evolving population. The analysis of the experiments shows that their feature selection method obtains results that are comparable with those obtained by considering all the features available. Thus, there is no performance increment, but only a reduction of the computational complexity.

Feature selection methods can be subdivided into two wide classes, filter and wrapper. Given a feature subset to be evaluated, filter functions take into account its statistical properties, while the wrapper ones use the performance achieved by a certain classifier trained on that subset. Filters methods generally involve a non-iterative computation on the dataset, which can be much faster than a classifier training session. In fact, implementing a classifier for evaluating the recognition rate attainable on a given subset, would require a costly training phase of such a classifier on a training set and a sample by sample labeling procedure on a test set. Moreover, filters methods evaluate intrinsic properties of the data, rather than the interactions of such data with a particular classifier: thus, the provided solutions should be more general, allowing good results to be obtained with a larger family of classifiers. The main drawback of filter methods is the fact that the objective function is generally monotonic, and this imply that the algorithm tends to select the full feature set as the optimal solution. This forces the user to select an arbitrary cut-off on the number of features to be selected. Wrapper methods generally achieve better recognition rates than filters ones since they are tuned tacking into account the specific interactions between the considered classifier and the dataset. These methods, however, are computationally expensive since they require that the learning procedure must be repeated for each feature subset, and the obtained results will be specific for the considered classifier.

Most of the approaches proposed in the context of handwriting recognition use wrapper methods (Chung et al., 1997; Kim et al., 2000; Oliveira et al., 2003; Nunes et al., 2004). Their main purpose is to reduce the number of features, keeping the recognition rate unchanged, or at most slightly worse.

Moving from these considerations, we propose a GA-based feature selection algorithm in which feature subsets are evaluated by means of a novel separability index. Our algorithm belongs to the filter method category and has been devised by extending the Fisher Linear Discriminant (Hart et al., 2001) method. Such method uses covariance matrices for estimating how the probability distributions of patterns are spread out in the considered N -dimensional space. Given a feature subset X , the Fisher's approach estimates the separability of the classes in X by taking into account two aspects: (i) how patterns belonging to a given class are spread out around the corresponding class mean vector (the centroid); (ii) distances among class mean vectors. Moreover, in order to compare subsets with different number of attributes, and to balance the effects of the monotonic trend of the separability index, we have added to the objective function a further term, suitable weighted, that takes into account the cardinality of the subspace to be evaluated. The proposed approach, thanks to the devised separability index, presents two main advantages: (i) it does not require that the dimensionality of the searched subspace (i.e. the actual number of features to be used) is a priori fixed; (ii) its performances are independent from the classification scheme.

The effectiveness of the proposed approach has been tested by using three standard databases of handwritten digits and a standard database of handwritten letters (uppercase and lowercase, totaling 52 classes), while the solutions found have been used to train different classifiers. The results have been compared with those obtained by using the whole feature set and with those obtained by using standard feature selection algorithms. The comparison outcomes confirmed the effectiveness of our approach.

The remainder of the paper is organized as follows: Section 2 discusses the feature selection problem, while Section 3 illustrates the proposed GA-based feature selection method. Section 4 describes the fitness function based on a specifically devised separability index used for subset evaluation. In Section 5 the experimental results are detailed, while some conclusions are eventually left to Section 6.

2. The feature selection problem

The goal of feature selection (FS) is that of reducing the number of features to be considered in the classification stage. This task is performed by removing irrelevant or noisy features from the whole set of the available ones. Feature selection is accomplished by reducing as much as possible the information loss due to the feature set reduction: thus, at list in principle, the selection process should not reduce classification performance. The feature selection process consists of three basic steps (see Fig. 1): a search procedure, a subset evaluation and a stopping criterion. A typical search procedure uses a search strategy for finding the optimal solution, according to a given subset evaluation criterion previously chosen. The search procedure is repeated until a stopping criterion is satisfied.

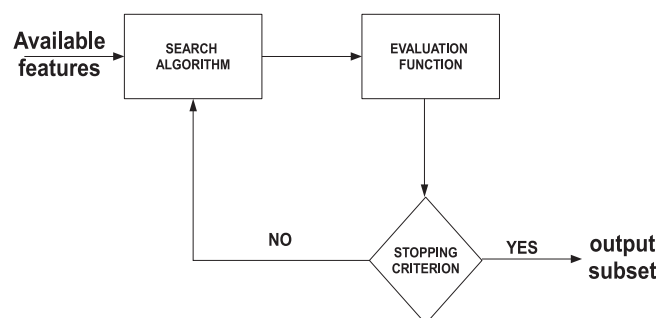


Fig. 1. The feature selection process.

Considering a generic application in which a set of samples (say \mathcal{Z}) must be classified, and assuming that the samples are represented by means of a set Y of N features, the feature selection problem can be formulated as follows: find the subset $X \subseteq Y$ of M features which optimizes an objective function J . Given a generic subset $X \subseteq Y$, $J(X)$ measures how well the patterns in \mathcal{Z} are discriminated by using the features subset X .

Example of statistical measures used by filter methods are the following: Distance (Ho and Basu, 2002), correlation (Guyon and Elisseeff, 2003; Hall, 2000), information (Ben-Bassat, 1982) and consistency (Dash and Liu, 2003). Distance-based criteria takes into account the geometrical characteristics of the class distributions in order to evaluate how well different classes are separated in the subset to be evaluated. Correlation measures use measures able to estimate the dependency between couple of variables. Such estimation can be used to find the correlation between a feature and a class. If the correlation between the feature x_1 and a given class c_i is higher than that between the feature x_2 and c_i , then the feature x_1 is preferred to x_2 for describing the class c_i . A slight variation of this criterion determines the dependence of a feature on the other ones; this value can be used to assess the redundancy degree of the features. Information measures, instead, evaluate the information gain from a given feature. The information gain of a feature x is defined as the difference between the a-priori uncertainty and the expected a-posteriori uncertainty of the class label given x ; the entropy measure can be used to estimate these uncertainties. Finally, the consistency measure of a feature subset, is determined counting the number of samples with the same feature values, but belonging to different classes.

Once the evaluation function $J(X)$ has been chosen, the feature selection problem becomes an optimization problem whose search space is the set of all the subsets of Y . As mentioned in the Introduction the size of this search space is exponential (2^N). As a consequence, the exhaustive search for the optimal solution becomes infeasible when a large number of features ($N > 50$) is involved. Search strategies like branch and bound (Yu and Yuan, 1993) have been proposed to strongly reduce the amount of evaluations, but the exponential complexity of the problem still remains. The exponential size of the search space for the feature selection problem makes appropriate the use of heuristic algorithms, for finding near-optimal solutions. Among these search algorithms, greedy search strategies are computationally advantageous but may lead to suboptimal solutions. They come in two flavors: forward selection and backward elimination. Forward selection strategies generate near-optimal feature subsets by a stepwise procedure which starts with an empty set. At each step the feature, among those not yet selected, that most increases the evaluation function J is added to the so far built subset; this procedure is repeated until a stop criterion is not satisfied. In backward elimination, instead, the whole subset of feature is initially considered, and at each step the feature that least reduce the evaluation function is eliminated. Both procedures are optimal at each step, but they cannot discover complex interactions among several features, as is the case in most of the real world feature selection problems. Then heuristic search algorithms, like genetic algorithms and simulated annealing (Meiri and Zahavi, 2006) seems to be appropriate for finding near-optimal solutions which take into account multiple interactions among several features.

3. The proposed method

In the framework of filter approaches to the feature selection problem, we propose a new method based on the use of genetic algorithms (GAs). These algorithms belong to the *evolutionary computation* paradigm (Goldberg, 1989), which has shown to be very

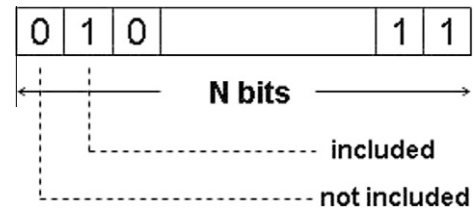


Fig. 2. An example of feature subset encoding by means of a bit string.

effective for solving optimization problems whose search spaces are high dimensional, discontinuous and complex. Mimicking the phenomena of natural evolution of species, GAs allows us to evolve a population of possible solutions, where each of them (denoted as an *individual* in the GA jargon) is represented as a binary string. Crossover and mutation operators are used to modify such strings in order to explore the search space, i.e. the set of all possible solutions. The method presented here has been implemented by using a generational GA, in which individuals are binary vectors each encoding a feature subset. More specifically, given a feature set Y having cardinality N , a subset X of Y ($X \subseteq Y$) is represented by an individual I having N elements whose i th element is set to 1 if the i th feature is included in X , 0 otherwise (see Fig. 2).

Besides the simplicity in the solution encoding, GAs are well suited for this class of problems because the search in this exponential space is very hard since interactions among features can be highly complex and strongly nonlinear. The algorithm starts by randomly generating a population of P individuals, whose values are set to 1 according a given probability (called *one_prob*). Such probability is usually set to low values (≈ 0.1) in order to force the early stage of the evolutionary search toward solutions having a small number of features. Then, the fitness of the generated individuals is evaluated by means of a suitably defined fitness function. This function takes into account how well the samples belonging to different classes are separated in the feature subset encoded by an individual, favoring at the same time the discovery of solutions containing a smaller number of features. After this evaluation phase, a new population is generated by first copying the best e individuals of the current population in order to implement an elitist strategy. Then $(P - e)/2$ couples of individuals are selected using the tournament method, which allows both loss of diversity and selection intensity to be controlled (Blickle and Thiele, 1996). The one point crossover operator is then applied to each of the selected couples, according to a given probability factor p_c . Afterwards, the mutation operator is applied. Then, the fitness function is computed according to the method illustrated in the next Section. Finally here individuals are added to the new population. The process just described is repeated for N_g generations. Note that it would be possible that some of the individuals generated according to the above process encode feature subset for which it is not possible to compute the fitness function. These solutions are simply discarded by the GA and a new offspring are generated by selecting other new parents in the current population.

4. Fitness function

The proposed fitness function takes into account two terms: in the first one, a function J , called separability index, measures the separability of the patterns belonging to different classes in the feature subset encoded by an individual. The second term takes into account the cardinality of the subset so as to favor solutions containing a smaller number of features.

The separability index J has been derived from the multiple discriminant analysis (MDA) approach. MDA is an extension to

C -class problems ($C > 2$) of the Fisher's Linear Discriminant (Hart et al., 2001), which has been defined for finding the best linear combination of features in case of two class problems.

In our case, assuming that each feature can be modeled as a random variable, the separability index J can be computed by using the covariance matrix of the whole feature set. Before providing the definition of J , let us recall some general properties of covariance matrix, Fisher's Linear Discriminant and multiple discriminant analysis.

Covariance matrix is the generalization of variance of a scalar variable to multiple dimensions. While variance measures the dispersion of the values of a random variable around its mean value, the covariance matrix of n variables measures how the joint probability distribution of the variables is spread out in the considered n -dimensional space around the mean vector. In particular, given n random variables $\{x_1, x_2, \dots, x_n\}$, each sampled by considering m values (stored in a $m \times n$ matrix D), the covariance matrix Σ is a $n \times n$ matrix in which the element at row i and column j represents the covariance between the variables x_i and x_j :

$$\Sigma[i, j] = \text{Cov}(x_i, x_j)$$

The covariance matrix is symmetric, in fact the generic element $\text{Cov}(x_i, x_j)$ is defined as follows:

$$\text{Cov}(x_i, x_j) = \frac{1}{m} \sum_{l=1}^m (D[l, i] - \mu_i)(D[l, j] - \mu_j) \quad (1)$$

where μ_i and μ_j are the mean values of the elements in the i th and j th column of D , respectively.¹

If samples belonging to different classes are represented in as points in a N -dimensional space, the k th class can be described by using its covariance matrix Σ_k , which is obtained by considering only the samples belonging to the class k . Such a matrix, in fact, reports information about the variability of k th class samples around their mean vector $\bar{\mu}_k$.

Denoting with c_0 and c_1 the classes to be discriminated and with ω the label associated to a sample \bar{x} , the Fisher's Linear Discriminant (FLD) (Hart et al., 2001) approach assumes that the conditional probability density functions $p(\bar{x}|\omega = c_0)$ and $p(\bar{x}|\omega = c_1)$ are normally distributed with mean and covariance parameters $(\bar{\mu}_0, \Sigma_0)$ and $(\bar{\mu}_1, \Sigma_1)$, respectively. It is also assumed that both classes have the same a priori probability ($P(c_0) = P(c_1)$). The goal of FLD approach is to find the vector \bar{w}^* representing the locus of the points where the samples are projected, which best separates the two classes. To this aim, the separation between the above distributions (denoted as class separability S in the following) has been computed as the ratio of the variance between the classes to the variance within the classes:

$$S(\bar{w}) = \frac{\sigma_{\text{between}}^2}{\sigma_{\text{within}}^2} = \frac{(\bar{w}^T \bar{\mu}_1 - \bar{w}^T \bar{\mu}_0)^2}{\bar{w}^T \Sigma_0 \bar{w} + \bar{w}^T \Sigma_1 \bar{w}} = \frac{\bar{w}^T (\bar{\mu}_1 - \bar{\mu}_0)(\bar{\mu}_1 - \bar{\mu}_0)^T \bar{w}}{\bar{w}^T (\Sigma_0 + \Sigma_1) \bar{w}} \quad (2)$$

According to the FDL theory, it can be shown that the best separation occurs when:

$$\bar{w} = \bar{w}^* = (\Sigma_0 + \Sigma_1)^{-1} (\bar{\mu}_1 - \bar{\mu}_0) \quad (3)$$

Note that the vector \bar{w}^* is normal to the discriminant hyperplane and can be computed through Eq. (3) only if the matrix resulting from the sum $(\Sigma_0 + \Sigma_1)$ is nonsingular, i.e. invertible.

When there are C classes to be discriminated, the analysis just described can be extended to find the $(C - 1)$ -dimensional subspace, which maximizes the class separability S . Note that if N is the number of available features, such subspace can be represented by means of a $N \times (C - 1)$ matrix $B\mathbf{W}$. Such a matrix is composed of

$(C - 1)$ N -dimensional vectors, representing the $(C - 1)$ projections on the transformed space. This approach is usually denoted in the literature as multiple discriminant analysis (MDA). In this case, the variances $\sigma_{\text{between}}^2$ and σ_{within}^2 can be expressed in terms of two matrices denoted as *within-class scatter matrix* Σ_W and *between-class scatter matrix* Σ_B :

$$\Sigma_W = \sum_{k=1}^C P(c_k) \Sigma_k$$

$$\Sigma_B = \sum_{k=1}^C P(c_k) (\bar{\mu}_k - \bar{\mu}_0)(\bar{\mu}_k - \bar{\mu}_0)^T$$

where $P(c_k)$ denotes the a priori probability of the k th class, Σ_k and $\bar{\mu}_k$ are the covariance matrix and the mean vector of k th class, respectively, and $\bar{\mu}_0$ denotes the overall mean:

$$\bar{\mu}_0 = \sum_{k=1}^C P(c_k) \bar{\mu}_k$$

Note that the *within-class scatter matrix* Σ_W measures the average spread of the classes about their mean vectors, while the *between-class scatter matrix* Σ_B measures the distances between each class mean vector and the overall mean.

According to the MDA theory, the class separability S can be measured as follows:

$$S(W) = \frac{|\mathbf{W}^T \Sigma_B \mathbf{W}|}{|\mathbf{W}^T \Sigma_W \mathbf{W}|} \quad (4)$$

where $|\cdot|$ indicates the determinant.

In this case the best separation is obtained by selecting the projections that give the best separation among classes. Such projections individuate the subspace represented by the $N \times (C - 1)$ matrix W^* , written as:

$$W^* = \Sigma_W^{-1} \Sigma_B$$

The matrix W^* allows a transformation of the original space in the projected space, thus a dimensionality reduction from a N to $C - 1$. Such a mapping is a good way to handle the curse of dimensionality but, as explained in (Hart et al., 2001), it can not possibly allow to obtain the minimum achievable error rate, especially in case of very large data set. Moreover the computational complexity of finding the optimal W^* is dominated by the calculation of the inverse of the *within-class scatter matrix* (Σ_W^{-1}). Note that the matrix Σ_W can be obtained by computing the covariances Σ_k ($k = 1, \dots, C$) through Eq. (1), but the computation of Σ_W^{-1} requires that $|\Sigma_W| \neq 0$ and this condition is not always verified. Finally, it is worth noticing that the multiple linear discriminant method does not directly represent a classification method, but rather it provides the subspace in which the classes are best separated: in this subspace a classification scheme must be defined in order to classify the patterns.

Moving from the above consideration, the basic idea of our approach is that of managing the curse of dimensionality by finding the optimal mapping from the original N -dimensional space to a M -dimensional one obtained by considering only a subset of M features among the whole set of N available ones. In other words, we have reformulated the feature extraction problem as a feature selection one, in which the effectiveness of the selected feature subspace is measured by using the class separability S defined in Eq. (4). Following this approach, the matrix \mathbf{W} is composed of M N -dimensional vectors, representing the axes of the transformed M -dimensional subspace. Such axes constitute the orthonormal basis of this M -dimensional subspace, and each of them coincides with one of the axes of the original N -dimensional feature space. Note that we do not assume that the selected features are independent, even if the considered subspaces are orthogonal: in fact, the

¹ Note that the element $\Sigma[i, i]$ represents the variance of the variable x_i .

orthogonality of the subspaces is simply a direct consequence of the basic assumption of our approach, which try to solve a feature selection problem and not a feature extraction one.

Let us now provide a formal definition of the separability index J : given an individual I , representing a feature subset X having dimensionality M , the separability index $J(I)$ is computed as follows:

$$J(I) = \text{tr} \left(\frac{\mathbf{W}^T \Sigma_B \mathbf{W}}{\mathbf{W}^T \Sigma_W \mathbf{W}} \right) \quad (5)$$

where \mathbf{W} is the transformation matrix² from the original N -dimensional space to the M -dimensional subspace corresponding to the subset X , while the symbol $\text{tr}(\cdot)$ is the trace operator. High values of the separability index $J(I)$ indicate that, in the subspace represented by the individual I , the centroids of the classes are well separated and, at the same time, the patterns are not too much spread out around their mean values. Without losing generality, we have modified the Eq. (4) using the trace operator: in fact, even if there is no mathematical equivalence between Eqs. (4) and (5), it has been demonstrated in (Fukunaga, 1990) that they provide the same set of optimal features. Therefore, we have chosen Eq. (5), which is computationally more effective.

Eventually, let us now define the fitness function used by our GA-based feature selection method: given an individual I , its fitness value $F(I)$ is computed by applying the formula:

$$F(I) = \frac{J(I)}{N} + K \frac{N - N_I}{N} \quad (6)$$

where N is the total number of features available, N_I is the cardinality of the subset represented by I (i.e. the number of bits equal to 1 in its chromosome) and K is a constant value used to weight the second term.

Note that in the first term, the separability index J has been divided by N so as to assure that both terms of Eq. (6) range from 0 to 1. The second term has been added since the first one exhibits a monotonic trend with the number of features. In fact, starting from a certain set of features, if we add any further feature, the separability index J do not decrease its value unless the new feature set makes the matrix Σ_W not invertible. In this case, obviously, the separability index cannot be computed and the new solution must be discarded. Thus, using as fitness only the first term of Eq. (6), the GA may produce solutions in which there are features not contributing to increase separability index J . These solutions are penalized by the presence of the second term. Finally, as regards the constant K , its role is to weight the second term in such a way that individuals having more features are favored only if they exhibits higher values for the separability index.

5. Experimental Results

In order to ascertain the effectiveness of the proposed approach, four real world datasets involving handwritten characters have been taken into account. Since our approach is stochastic, as well

² For example, if $N = 4$ is the cardinality of the whole feature space and the individual $I = (0, 1, 0, 1)$ encodes the subspace X having cardinality 2 (features 2 and 4) the matrix \mathbf{W} is the following:

$$\mathbf{W} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{pmatrix}$$

Note that if Σ is the 4×4 covariance matrix, the product $\mathbf{W}^T \Sigma \mathbf{W}$ gives as result the 2×2 matrix Σ' , which represents the projection of Σ on the subspace encoded by the individual I .

as all the EC-based algorithms, 20 runs have been performed for each experiment carried out. The ability of our system in finding good subsets has been evaluated by measuring the performance obtained with three classification algorithms: support vector machine (SVM) (Vapnik, 1998), multiple layer perceptron (MLP) and k -Nearest Neighbor (k -NN) (Hart et al., 2001). To this purpose, at the end of each run, the best feature subset found has been used for training the considered classifiers and for evaluating their performance.

As regards the SVM, we used the implementation provided by the LIBSVM public domain software (Chang and Lin, 2001). In the experiments described in the following we used three different kernels: Radial Basis Function (RBF), Polynomial and Sigmoidal. For the k -NN and MLP classification algorithms, we used the implementation provided by the WEKA tool (Hall et al., 2009). As regards the evolutionary parameters, shown in Table 1, they have been heuristically found performing a set of preliminary trials; this set of parameters has been used for all the experiments reported below.

The main objectives of the experiments are the following:

- To investigate the influence of the cost factor K on the performance of the system.
- To study if the selected feature subsets are independent of the choice of the used classifier. To this aim, the subsets found have been used to train the considered classification algorithms with different sets of parameters.
- To perform a comparison of the obtained results with those of other feature selection algorithms reported in the literature.

In the following, the datasets used and the experiments performed are detailed.

5.1. The datasets

We have used in the experiments three standard databases of handwritten digits and a standard database of handwritten letters (uppercase and lowercase, totaling 52 classes). Two of them (OPTODIGIT and MFEAT) are publicly available from the UCI machine learning repository (Frank and Asuncion, 2010), the third one (MNIST) has been made available by the New York University (LeCun and Cortes, 2010), while the last one (NIST-SD19) is provided by the National Institute of Standard Technologies (Grother, 1995).

The MFEAT dataset (multiple features dataset) contains 2000 instances of handwritten digits, 200 for each digit, extracted from a collection of Dutch utility maps. Data are described by using six different sets of features, totaling 649 features. Each set of features has been used to describe all the handwritten digits, and arranged in separate datasets. This implies that we have six datasets (DS1, ..., DS6), each containing 2000 samples. For each dataset, the type of features and their number are the following:

- DS1: 76 Fourier coefficients of the character shapes.
DS2: 47 Zernike moments.

Table 1

The evolutionary parameters.

| Parameter | symbol | value |
|-----------------------|--------|-------|
| population size | P | 100 |
| tournament size | t | 5 |
| elitism size | e | 1 |
| crossover probability | p_c | 0.6 |
| mutation probability | p_m | $1/N$ |
| generation number | N_g | 1000 |

DS3: 6 morphological features.

DS4: 64 Karhunen-Love coefficients.

DS5: 240 pixel averages in 2×3 windows.

DS6: 216 profile correlations.

More details about the feature sets can be found in (Breukelen et al., 1997). Starting from the provided datasets we generated a further dataset (DS) obtained by merging all the descriptions included in the previous ones, in such a way to describe each sample by the whole set of 649 available features. From the generated dataset DS, 70 samples per class have randomly extracted to build a training set (TR). The remaining data have been used to build a test set (TS) including 130 samples per class. Summarizing, TR contains 700 samples, while TS contains 1300 samples.

The second considered dataset is the optical recognition of handwritten digits dataset (OPTODIGIT). It contains 5620 samples equally distributed among the ten classes. Each sample is described by 64 features. Such data have been obtained by preprinted forms, extracting normalized 32×32 bitmaps of handwritten digits. Each bitmap is divided into non-overlapping blocks of 4×4 and the number of black pixels are counted in each block. This generates an input matrix of 8×8 where each element is an integer in the range $[0, 16]$. As a consequence, a character is represented by a feature vector of 64 elements where each element contains a value of a 8×8 matrix. In order to build a training set, 3820 samples have been randomly picked up from the original dataset, while the remaining ones (1800) have been used as test set.

The third dataset taken into account is MNIST. Such dataset was constructed from NIST's Special Database 3 and Special Database 1 which contain binary images of handwritten digits (LeCun and Cortes, 2010). The original black and white images from NIST were size normalized to fit in a 20×20 pixel box while preserving their aspect ratio. The resulting images contain gray levels as a result of the anti-aliasing technique used by the normalization algorithm. The images were centered in a 28×28 image by computing the center of mass of the pixels, and translating the image so as to position this point at the center of the 28×28 field. Finally, the MNIST training set is composed of 60,000 samples, while the test set contains 10,000 samples. The total number of features used to describe samples is 780, even if some of them assume values different from zero in less than 1% of samples. For this reason,

we discarded such features considering in our experiment only 485 features.

Finally, the NIST-SD19 database (NIST in the following), contains binary alphanumeric characters. In particular, we have considered handwritten letters, uppercase and lowercase, corresponding to 52 classes. The handwriting sample form *hsf4*, containing 23,941 characters, has been used as training set, while the handwriting sample form *hsf7*, containing 23,670 characters, has been used as test set. *hsf4* has 11,941 uppercase characters and 12,000 lowercase ones, while *hsf7* has 12,092 uppercase characters and 11,578 lowercase ones. In each form, characters are segmented and stored in 128×128 pixel images, each associated to one of the 52 classes to be discriminated. Samples are represented by using the features proposed in (Oliveira et al., 2002). Each character is described by a feature vector containing the measures associated to different parts of the image. More specifically, the image is divided in six parts and, for each part, 22 features are computed, totaling 132 features.

5.2. The Constant K

Several experiments have been performed for analyzing how the constant K affects the behavior of our feature selection method, in terms of both number of features and classification performance. Such a performance refers to the use of RBF SVM classifiers. As expected, the higher K the lower the number of selected features and the obtainable classification results. For the sake of clarity, we have not reported in the following plots the values of K for which the performance is too much degraded.

In the Figs. 3–6 we have reported both the recognition rate (RR) and the number of features (N_i/N) for the datasets MFEAT, OPTODIGIT, MNIST and NIST, respectively. For each value of K we have reported the average result over the 20 performed runs. Note that we have shown in the figures only the standard deviations, which do not assume negligible values.

As regards the MFEAT dataset, the Fig. 3 shows that for a large interval of K values, the performance is almost constant obtaining the highest value for $K = 0.01$. This value has been chosen hereafter in the experiments for MFEAT dataset, even if values up to 1.0 may be used slightly reducing classification performance, but strongly reducing also the dimension of the feature space. Similar

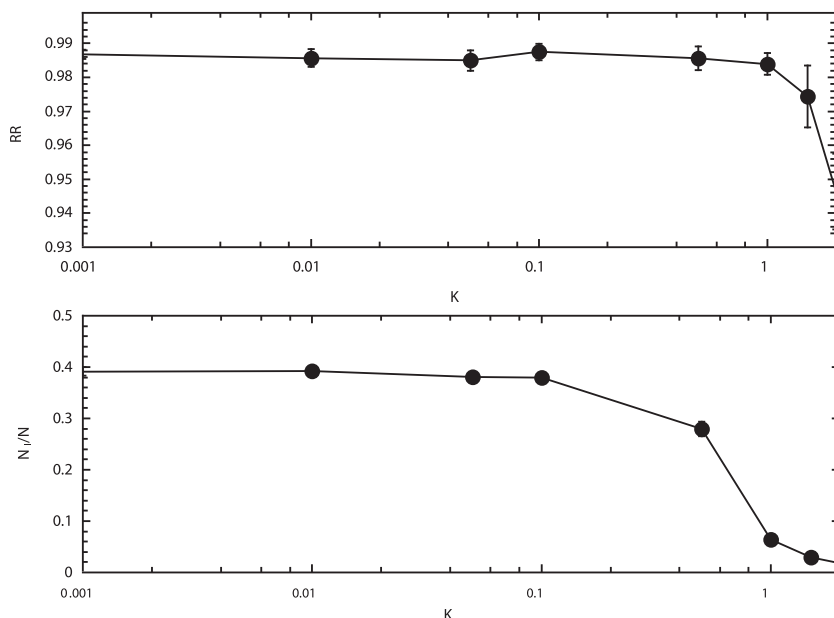


Fig. 3. Recognition rate RR (top) and number of features N_i (bottom) as a function of the constant K for MFEAT dataset.

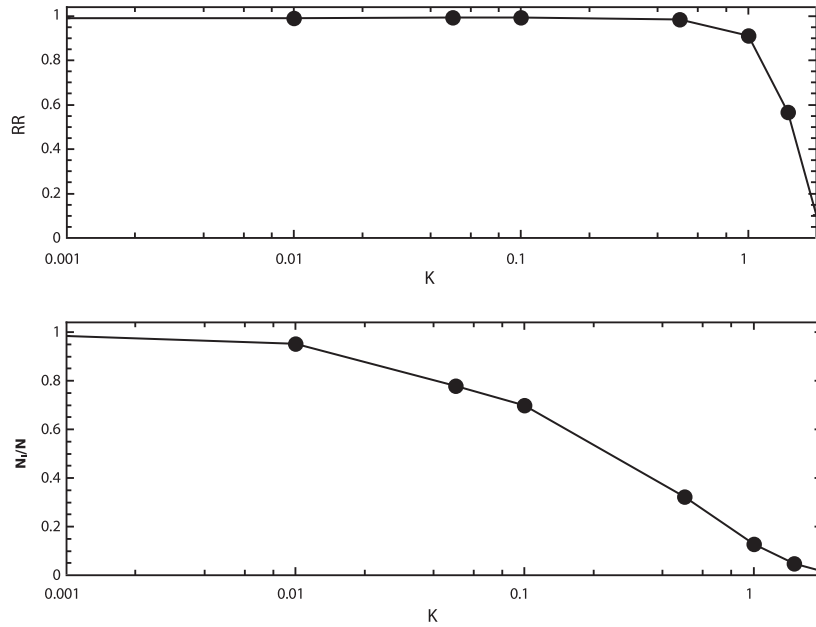


Fig. 4. Recognition rate RR (top) and number of features N_i (bottom) as a function of the constant K for OPTODIGIT dataset.

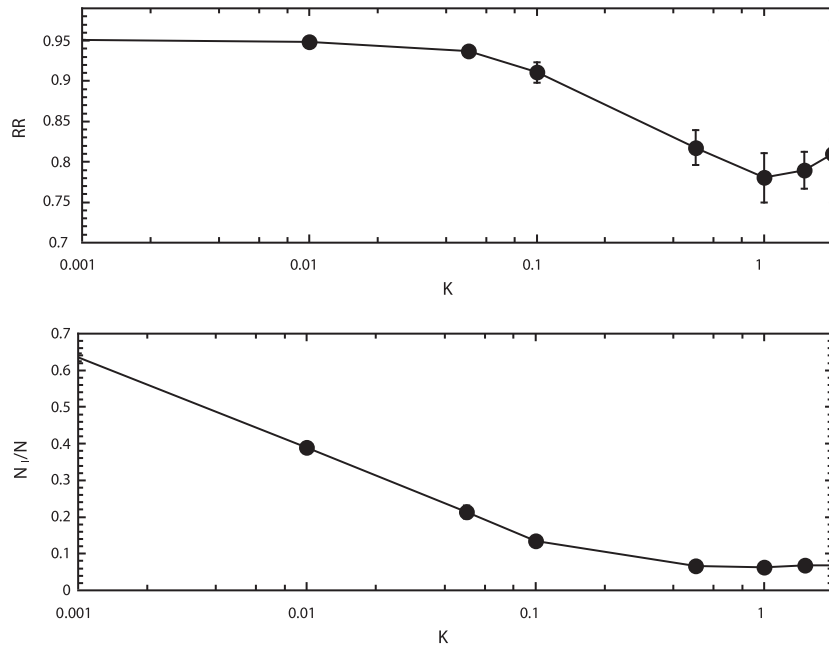


Fig. 5. Recognition rate RR (top) and number of features N_i (bottom) as a function of the constant K for MNIST dataset.

considerations hold for the other two datasets: in particular for OPTODIGIT dataset the optimal value chosen for K is 0.05, even if the value $K = 0.5$ results in a very small reduction of the recognition rate, but using only less than 30% of the available features. For the MNIST dataset, the optimal value chosen for K is 0.001, even if the value $K = 0.5$ corresponds to a small reduction of the recognition rate, but using only the 10% of the available features. Finally, as regards NIST, the optimal value chosen for K is 0.1, even if the value $K = 0.1$ corresponds to a small reduction of the recognition rate but using only less than 20% of the available features.

Summarizing, as shown in Figs. 3–6, the trend of RR as a function of K is very regular: there is an interval of values of K in which RR assumes an almost constant or slightly increasing trend and

then it rapidly decreases. Thus, exploiting the regularity of RR curve, a general experimental procedure for setting the value of the parameter K , is the following: increase the values of K until RR exhibits slight variations; as soon as RR starts to rapidly decrease for a certain value of K , select the previous K value as the optimal one.

5.3. Behavior analysis

Since our feature selection method does depend on any specific classifiers, we want to evaluate both the generality and the effectiveness of our results by using different classification algorithms. As a consequence, for each dataset, the best feature subset provided by the GA has been tested by using different classifiers

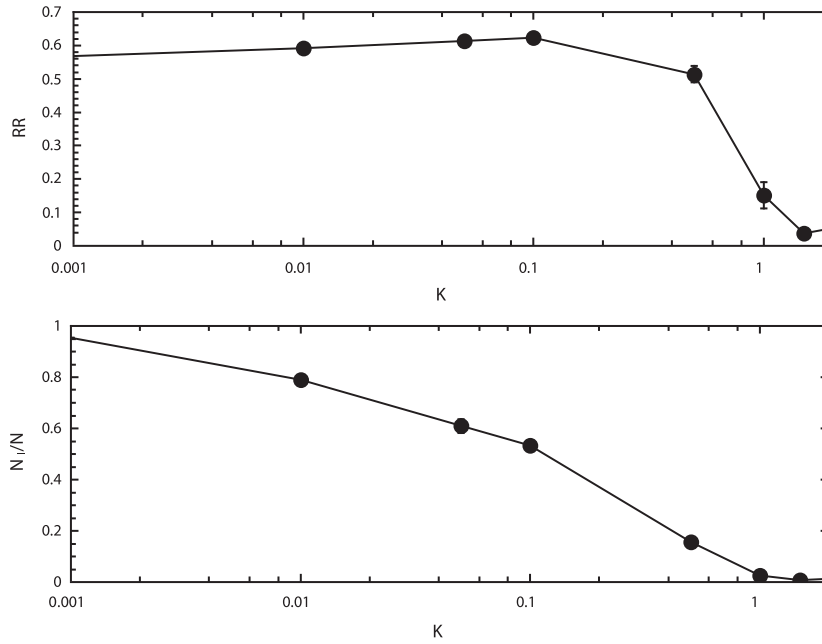


Fig. 6. Recognition rate RR (top) and number of features N_i (bottom) as a function of the constant K for NIST dataset.

obtained varying the configuration parameters of the three considered classification schemes (SVM, MLP and k -NN). The obtained performances have been compared with those achieved by using the whole feature sets.

The purpose of such a comparison is twofold: on the one hand, we want to verify if the selected features allows us to improve the performance with respect to those relative to all the available features. On the other hand, we want to understand if the differences

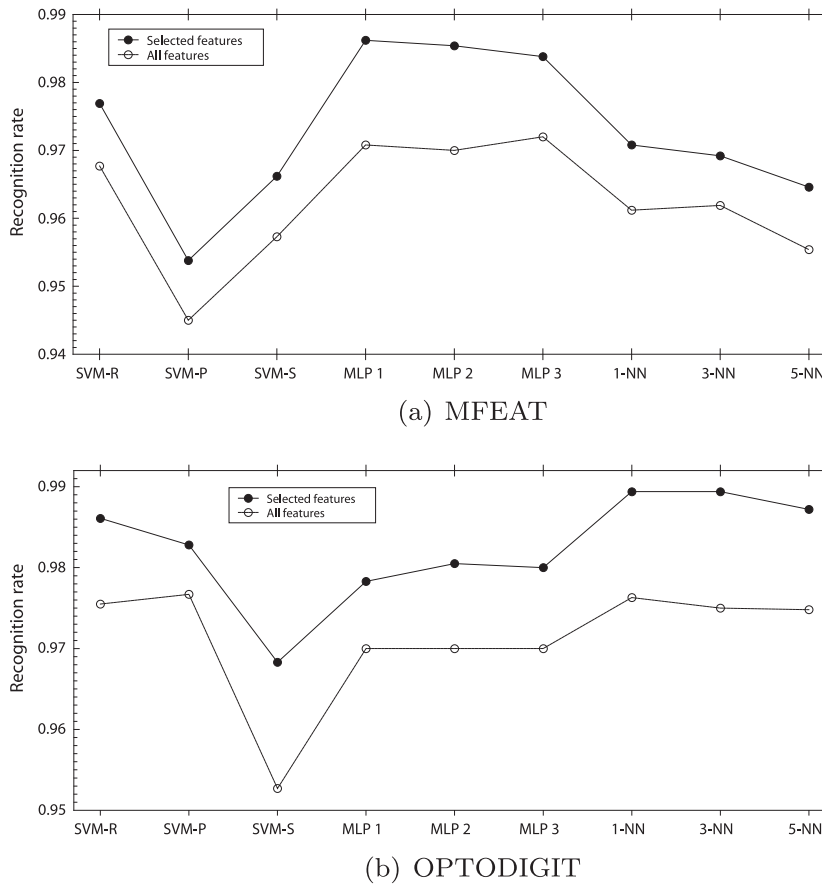


Fig. 7. Behavior analysis for the datasets MFEAT and OPTODIGIT.

among the considered classifiers actually depend on the selected features.

For each of the three algorithm, three different configuration parameter values have been considered, totaling nine classifiers. More specifically, for the SVM three different kernels have been used: radial basis function (RBF), Polynomial (with the degree set to three), and the sigmoid. For the MLP, different classifiers have been obtained by changing the number of neurons in the hidden layer. The following values have been used: 50, 100 and 200. Finally, for the k -NN, the following k values have been considered: 1, 3, 5. The performances of the 9 classifiers have been compared by means of the 10-fold cross validation method. The results are shown in Figs. 7 and 8. The acronyms SVM-R, SVM-P and SVM-S, stand for RBF SVM, Polynomial SVM and sigmoidal SVM, respectively. As regards the acronyms MLP-1, MLP-2 and MLP-3 they refer to the MLPs with 50, 100 and 200 hidden neurons, respectively. The acronyms of the k -NN are self-explanatory.

The figure shows that the feature subsets selected by our method always give better results than those obtained by using the whole set of features. Moreover, the trend of the performance obtained by using of the selected feature subsets is very similar to that of the whole feature set, confirming the generality of the proposed feature selection method.

5.4. Comparison findings

In order to test the effectiveness of the proposed system, our results have been compared with those obtained by four widely used feature selection techniques. Such techniques combines a search strategies and a subset evaluation criterion. Moreover, the

performance obtained using the whole feature set has been also considered.

As regards the search strategies we used the Best First (BF) (Xu et al., 1988) and the Linear Forward (LF) (Gütlein et al., 2009) ones. The former strategy searches solutions by using a greedy hill-climbing technique. It starts with the empty set of features and adds new features according to the best first search algorithm (Pearl, 1984). The latter one, instead, represents an extension of the Best First strategy. Such technique reduces the number of attribute expansions in each forward selection step. It is faster than BF and it generally finds smaller subsets. Our search strategy, based on the use of a GA, exhibits an higher computational cost, since it requires a number of operations equal to the number of individuals in the population by the number of generations.

As subset evaluation criteria we have considered the following ones: Feature-Class Correlation, Consistency Criterion and three different wrapper evaluation functions.

The Feature-Class Correlation (FCC) (Hall et al., 1998) evaluates a feature subset by measuring the correlation among its features and the classes: it prefers subsets highly correlated with classes, but having low correlation among features. The Consistency Criterion (CC) (Liu and Setiono, 1996) evaluates the worth of the feature subsets by using a consistency index measuring how well samples belonging to different classes are separated. As concerns the wrapper functions we used the same three classifiers used for the performance evaluation, namely RBF SVM, a MLP with 200 hidden neurons and the 3-NN.

The computational complexity of our evaluation function is lower than those exhibited by both FCC and CC. In fact, we have used the training set data only for computing, once and forever, the within-class and the between-class scatter matrices Σ_W and

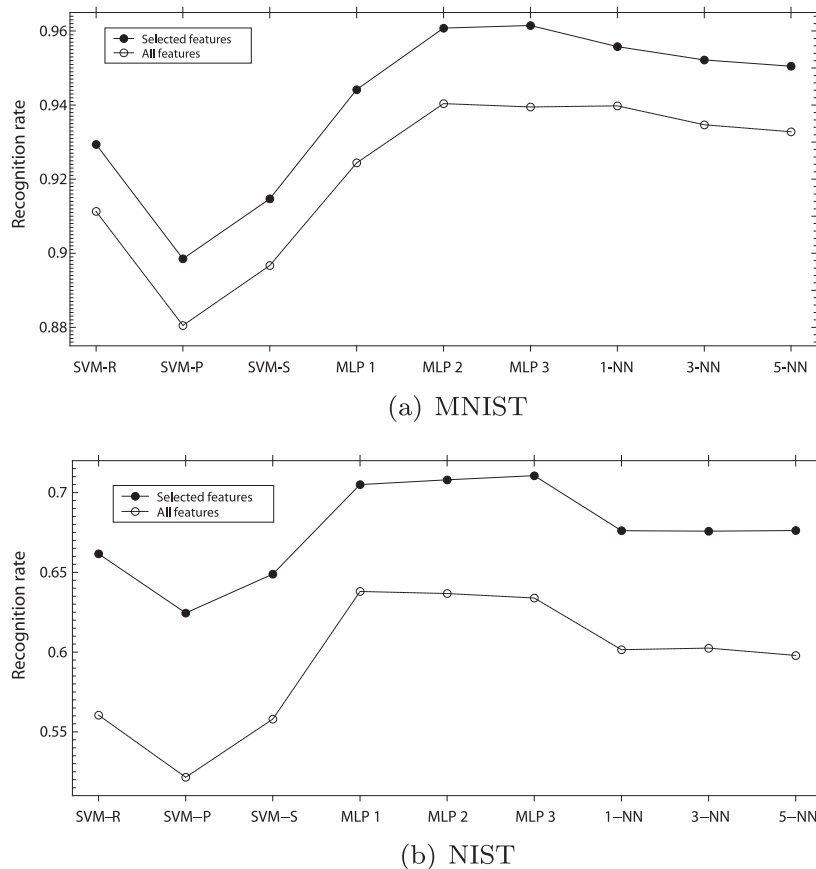


Fig. 8. Behavior analysis for the datasets MNIST and NIST.

Σ_B . The other evaluation functions, FCC and CC, on the contrary, compute the feature correlation and the consistency index, respectively, for each subset to be evaluated. Finally, the computational complexity of the wrapper functions is the highest one, since they require to perform, for each subset to be evaluated, the training of the classifier and the evaluation of the corresponding recognition rate.

The results are shown in Table 2. In the table, the methods considered for the comparison are denoted by acronyms. For each acronym the first part denotes the evaluation function, while the second one the search strategy. For instance, FCC-BF denotes a comparing method using the Feature–Class Correlation as evaluation function and the Best First search strategy.

Note that the accuracy results have been obtained by using the 10-fold cross validation. For each dataset the number of selected features are also reported. The last row shows, for each dataset, the average values of the performance differences between the proposed approach and the methods taken into account for the comparison.

The data reported in the table shows that our approach achieves better results than those obtained by the comparing methods. The last row of this table, reports the mean values of such improvements, computed averaging the differences between our method and the others for each dataset. It can be seen that such improvements varies from about 1.7% to 5.6%. As regards the number of selected features, our approach selects more features than the compared ones. This is due to the fact that we chose the value of the constant K , which maximizes the recognition rate even if larger

values of K allowed a slight performance reduction, but a considerable smaller number of features.

In order to have a fair comparison with methods selecting small numbers of features, and to test the performance of our system when as less as possible features should be used, we performed a further set of experiments. More specifically, for each dataset, we chose in the fitness function a value for the parameter K , such as to obtain a number of features similar to that provided by the methods selecting the smallest number of features (CC-LF and CC-BF). The results of this comparison are shown in Table 3. Also in this case, accuracy results have been obtained by using the 10-fold cross validation.

The data in the table show that even when high values of K are used (strong feature reduction), our method performs almost always better than the other considered ones. In particular, as concerns MFEAT and OPTODIGIT datasets, our approach selects the same number of feature, or a slightly higher one, but achieves better results. The performance increments vary from 1% (MFEAT, KNN classifier) to 5% (OPTODIGIT, MLP classifier). As for MNIST dataset, our method selects a number of features higher than the that of CC-BF method (34 vs 13), but in this case the performance improvements are more relevant, ranging from 8.8% (SVM) to 16.1% (KNN). Finally, as regards the NIST dataset, the results of our method are slightly worst than those of CC-LF method (approximately from 1% to 2% less). It should be considered, however, that these are the less favorable conditions for our method: in fact, as the value of K grows, the second term in the fitness function becomes more and more relevant with respect to the first one (see

Table 2
Comparison results.

| | Cl. | MFEAT | | MNIST | | OPTODIGIT | | NIST | |
|-------------------|------------|-------|-----|-------|-----|-----------|----|-------|-----|
| | | Acc. | # | Acc. | # | Acc. | # | Acc. | # |
| All | SVM | 96.77 | 649 | 91.13 | 489 | 97.55 | 64 | 56.05 | 132 |
| | MLP | 97.2 | | 93.95 | | 97 | | 63.39 | |
| | KNN | 96.19 | | 93.47 | | 97.5 | | 60.25 | |
| Our Method | SVM | 97.69 | 233 | 92.94 | 272 | 98.61 | 44 | 66.16 | 61 |
| | MLP | 98.38 | | 96.15 | | 98 | | 71.06 | |
| | KNN | 96.92 | | 95.22 | | 98.94 | | 67.58 | |
| FCC-BF | SVM | 97.08 | 133 | 91.76 | 144 | 97.55 | 35 | 62.1 | 64 |
| | MLP | 97.04 | | 94.16 | | 97.07 | | 64.26 | |
| | KNN | 96.35 | | 93.92 | | 97.72 | | 63.58 | |
| FCC-LF | SVM | 97.15 | 106 | 88.36 | 58 | 97.55 | 35 | 62.21 | 57 |
| | MLP | 97.54 | | 88.23 | | 97.07 | | 64.57 | |
| | KNN | 96.77 | | 88.85 | | 97.72 | | 64.08 | |
| CC-BF | SVM | 93.46 | 6 | 73.65 | 13 | 89.32 | 9 | 58.55 | 17 |
| | MLP | 91.92 | | 68.09 | | 83.81 | | 64.39 | |
| | KNN | 93.08 | | 70.23 | | 87.54 | | 62.25 | |
| CC-LF | SVM | 93.92 | 6 | 77 | 20 | 90.43 | 9 | 58.95 | 14 |
| | MLP | 93.08 | | 72.79 | | 83.92 | | 64.89 | |
| | KNN | 93.92 | | 78.23 | | 88.65 | | 62.65 | |
| MLP-BF | SVM | 96.92 | 25 | 91.57 | 65 | 96.67 | 33 | 59.24 | 48 |
| | MLP | 96.98 | | 94.74 | | 97.07 | | 64.89 | |
| | KNN | 95.94 | | 94 | | 96.56 | | 64.11 | |
| MLP-LF | SVM | 97.19 | 14 | 91.68 | 43 | 96.85 | 26 | 58.48 | 41 |
| | MLP | 97.28 | | 94.85 | | 97.1 | | 65.07 | |
| | KNN | 96.02 | | 94.15 | | 96.88 | | 64.48 | |
| SVM-BF | SVM | 96.94 | 29 | 91.76 | 68 | 97.47 | 35 | 61.2 | 46 |
| | MLP | 96.75 | | 93.16 | | 97.04 | | 66.23 | |
| | KNN | 95.94 | | 93.14 | | 96.63 | | 65.25 | |
| SVM-LF | SVM | 97 | 13 | 91.83 | 47 | 97.67 | 30 | 61.4 | 40 |
| | MLP | 96.92 | | 93.27 | | 97.44 | | 66.45 | |
| | KNN | 96.01 | | 93.18 | | 97.63 | | 65.55 | |
| KNN-BF | SVM | 95.57 | 31 | 91.7 | 71 | 97.14 | 39 | 58.48 | 59 |
| | MLP | 95.3 | | 94.2 | | 96.65 | | 65.07 | |
| | KNN | 95.88 | | 94.33 | | 97.5 | | 64.48 | |
| KNN-LF | SVM | 95.57 | 17 | 91.85 | 41 | 97.6 | 32 | 58.68 | 48 |
| | MLP | 97.3 | | 94.27 | | 96.83 | | 65.7 | |
| | KNN | 95.88 | | 94.36 | | 97.91 | | 65.08 | |
| Avg diff. | | 1.7 | | 5.6 | | 3.1 | | 5.48 | |

Table 3
Further comparison results.

| Datasets | Cl. | Our method | | Others | | |
|-----------|-----|------------|----|--------|-------|----|
| | | Acc. | # | Method | Acc. | # |
| MFEAT | SVM | 95.0 | 6 | CC-LF | 93.92 | 6 |
| | MLP | 95.07 | | | 93.08 | |
| | KNN | 94.85 | | | 93.92 | |
| MNIST | SVM | 82.5 | 34 | CC-BF | 73.65 | 13 |
| | MLP | 83.49 | | | 68.09 | |
| | KNN | 86.26 | | | 70.23 | |
| OPTODIGIT | SVM | 92.82 | 12 | CC-LF | 90.43 | 9 |
| | MLP | 88.87 | | | 83.92 | |
| | KNN | 92.43 | | | 88.65 | |
| NIST | SVM | 57.88 | 26 | CC-LF | 58.95 | 14 |
| | MLP | 62.57 | | | 64.89 | |
| | KNN | 61.45 | | | 62.65 | |

Eq. (6)), and the evolutionary algorithm tends to favor solutions using a small number of features even if they exhibit a low value of the separability index. This behavior has been analyzed in Subsection 5.2 discussing the set of experiments for finding the optimal value of K .

6. Conclusions

In the framework of handwriting recognition problems, we have presented a feature selection method for improving classification performance. The devised approach uses a GA-based feature selection algorithm for detecting feature subsets where the samples belonging to different classes are well discriminated. The proposed method does not require that the dimensionality of the searched subspace is a priori fixed. Candidate feature subsets are evaluated by means of a novel evaluation function based on the Fisher linear discriminant. Such evaluation function uses covariance matrices for estimating how the probability distributions of patterns are spread out in the considered representation space. Moreover, in order to balance the effects of the monotonic trend of the evaluation function, a further term, suitable weighted, has been added which takes into account the cardinality of the subspace to be evaluated. The experiments have been performed by using four standard databases and the solutions found have been tested with different classification methods. The results have been compared with those obtained by using the whole feature set and with those obtained by using standard feature selection algorithms. The comparison outcomes confirmed the effectiveness of our approach.

From the experimental results we can also draw the following observations:

1. The choice of an appropriate value for the constant K is not critical. In fact, the trend of the recognition rate as a function of K is nearly constant for a wide range of K values and, starting from a certain point, it rapidly decreases. This behavior allows us to easily identify an appropriate range of values for K . More specifically, in application for which it is mandatory to maximize performance, the value of K corresponding to the highest recognition rate will be chosen. On the contrary, if we can accept a slight performance decrease, values of K corresponding to a lower number of features may be chosen.
2. The accuracy gain obtained by our method with respect to the compared methods justifies the computational time required by the feature selection algorithm proposed.
3. The separability index assumes that class distributions are normal. Even if this assumption is difficult to verify for real data, it represents a reasonable and a widely used approximation. Moreover, the results obtained on different datasets seems to

prove that this assumption does not negatively affect the obtainable performance in handwriting recognition applications.

Finally, it is worth noting that the proposed approach also shows the following interesting properties: (i) it is independent of the classification system used (ii) its computational time is independent of the training set size.

The second property derive from the fact that we have used the training set only for computing, once and forever, the covariance matrices Σ_k ($k = 1, \dots, C$), the mean vectors $\bar{\mu}_k$ ($k = 1, \dots, C$) and the overall mean vector $\bar{\mu}_0$. These information are needed for computing the within-class and the between-class scatter matrices Σ_W and Σ_B . This property is another consequence of the adopted filter approach, which makes our method suitable when large datasets are taken into account.

References

- Ahlgren, R., Ryan, H., Swonger, C., 1971. A character recognition application of an iterative procedure for feature selection. *IEEE Trans. Comput.* C-20 (9), 1067–1075.
- Ben-Bassat, M., 1982. Pattern recognition and reduction of dimensionality. In: Krishnaiah, P., Kanal, L. (Eds.), *Handbook of Statistics-II*. North Holland, 1982, pp. 773–791.
- Blickle, T., Thiele, L., 1996. A comparison of selection schemes used in evolutionary algorithms. *Evol. Comput.* 4 (4), 361–394.
- Breukelen, M.V., Duin, R., Tax, D., den Hartog, J., 1997. Combining classifiers for the recognition of handwritten digits. In: 1st internat. workshop on statistical techniques in, pattern recognition, pp. 13–18.
- Chang, C.-C., Lin, C.-J., 2001. LIBSVM: A library for support vector machines, available at <<http://www.csie.ntu.edu.tw/~cjlin/libsvm>>.
- Chouaib, H., Terrades, O.R., Tabbone, S., Cloppet, F., Vincent, N., 2008. Feature selection combining genetic algorithm and adaboost classifiers. In: *Proc. of Internat. Conf. on Pattern Recognition (ICPR 2008)*. IEEE.
- Chung, K., Yoon, J., 1997. Performance comparison of several feature selection methods based on node pruning in handwritten character recognition. In: *Proc. 4th Internat. Conf. on Document Analysis and Recognition (ICDAR 97)*. IEEE Computer Society, pp. 11–15.
- Cordella, L., De Stefano, C., Fontanella, F., Marrocco, C., 2008. A feature selection algorithm for handwritten character recognition. In: 19th Internat. Conf. on Pattern Recognition (ICPR 2008), 2008, pp. 1–4.
- Cordella, L., De Stefano, C., Fontanella, F., Marrocco, C., Scotto di Freca, A., 2010. Combining single class features for improving performance of a two stage classifier. In: 20th Internat. Conf. on Pattern Recognition (ICPR 2010), pp. 4352–4355.
- Dash, M., Liu, H., 2003. Consistency-based search in feature selection. *Artif. Intell.* 151 (1–2), 155–176.
- De Stefano, C., Fontanella, F., Marrocco, C., Schirinz, G., 2007. A feature selection algorithm for class discrimination improvement. In: *Geoscience and Remote Sensing, Symposium, 2007 (IGARSS07)*, pp. 425–428.
- Fodor, I., 2002. A survey of dimension reduction techniques, Tech. rep., Center for Applied Scientific Computing, Lawrence Livermore National Laboratory.
- Frank, A., Asuncion, A., 2010. UCI machine learning repository. URL <<http://archive.ics.uci.edu/ml>>
- Fukunaga, K., 1990. *Introduction to Statistical Pattern Recognition*, second ed. Academic Press Professional, Inc., San Diego, CA, USA.
- Goldberg, D.E., 1989. *Genetic Algorithms in Search Optimization and Machine Learning*. Addison-Wesley.
- Grother, P.J., 1995. Nist special database 19. URL <<http://www.nist.gov/srd/nistsd19.cfm>>
- Gütlein, M., Frank, E., Hall, M., Karwath, A., 2009. Large-scale attribute selection using wrappers. In: *Proc. IEEE Symposium on Computational Intelligence and Data Mining*. IEEE, pp. 332–339.
- Guyon, I., Elisseeff, A., 2003. An introduction to variable and feature selection. *J. Machine Learn. Res.* 3, 1157–1182.
- Hall, M.A., 2000. Correlation-based feature selection for discrete and numeric class machine learning. In: *ICML '00: Proc. Seventeenth Internat. Conf. on Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 359–366.
- Hall, M.A., 1998. Correlation-based feature subset selection for machine learning, Ph.D. thesis, University of Waikato, Hamilton, New Zealand.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I.H., 2009. The weka data mining software: An update. *SIGKDD Explor.* 11 (1), 10–18.
- Hart, P.E., Duda, R.O., Stork, D.G., 2001. *Pattern Classification*. John Wiley & sons Inc.
- Ho, T.K., Basu, M., 2002. Complexity measures of supervised classification problems. *IEEE Trans. Pattern Anal. Machine Intell.* 24 (3), 289–300.
- Kim, G., Govindaraju, V., 1997. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Trans. Pattern Anal. Machine Intell.* 19 (4), 366–379.

- Kim, G., Kim, S., 2000. Feature selection using genetic algorithms for handwritten character recognition. In: Proc. 7th Internat. Workshop on Frontiers in Handwriting Recognition (IWFHR 2000). Amsterdam, 2000, pp. 103–112.
- Kudo, M., Sklansky, J., 2000. Comparison of algorithms that select features for pattern recognition. *Pattern Recognition* 33 (1), 25–41.
- Kwak, N., Choi, C.-H., 2002. Input feature selection for classification problems. *IEEE Trans. Neural Networks* 13 (1), 143–159.
- LeCun, Y., Cortes, C., 2010. MNIST handwritten digit database. URL <<http://yann.lecun.com/exdb/mnist>>
- Liu, H., Setiono, R., 1996. A probabilistic approach to feature selection – a filter solution. In: 13th Internat. Conf. on Machine Learning, pp. 319–327.
- Meiri, R., Zahavi, J., 2006. Using simulated annealing to optimize the feature selection problem in marketing applications. *Eur. J. Oper. Res.* 171 (3), 842–858.
- Nunes, C.M., Britto Jr., A.d.S., Kaestner, C.A.A., Sabourin, R., 2004. An optimized hill climbing algorithm for feature subset selection: Evaluation on handwritten character recognition. In: Proc. 9th Internat. Workshop on Frontiers in Handwriting Recognition (IWFHR'04). IEEE Computer Society, Washington, DC, USA, pp. 365–370.
- Oh, I.-S., Lee, J.-S., Moon, B.-R., 2004. Hybrid genetic algorithms for feature selection. *IEEE Trans. Pattern Anal. Machine Intell.* 26 (11), 1424–1437.
- Oliveira, L.L.S., Sabourin, R., Bortolozzi, F., Suen, C., 2002. Automatic recognition of handwritten numerical strings: A recognition and verification strategy. *IEEE Trans. Pattern Anal. Machine Intell.* 24 (11), 1438–1454.
- Oliveira, L.S., Sabourin, R., Bortolozzi, F., Suen, C., 2003. A methodology for feature selection using multi-objective genetic algorithms for handwritten digit string recognition. *Internat. J. Pattern Recognit. Artif. Intell.* 17, 2003.
- Pearl, J., 1984. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Addison-Wesley, Reading, MA.
- Pudil, P., Novovičová, J., Kittler, J., 1994. Floating search methods in feature selection. *Pattern Recognition Lett.* 15 (11), 1119–1125.
- Shi, D., Shu, W., Liu, H., 1998. Feature selection for handwritten chinese character recognition based on genetic algorithms. In: 1998 IEEE Internat. Conf. on Systems, Man, and Cybernetics, vol. 5, pp. 4201–4206.
- Siedlecki, W., Sklansky, J., 1989. A note on genetic algorithms for large-scale feature selection. *Pattern Recognition Lett.* 10 (5), 335–347.
- Somol, P., Pudil, P., Kittler, J., 2004. Fast branch and bound algorithms for optimal feature selection. *IEEE Trans. Pattern Anal. Machine Intell.* 26 (7), 900–912.
- Vapnik, V.N., 1998. *Statistical Learning Theory*. Wiley-Interscience.
- Xu, L., Yan, P., Chang, T., 1988. Best first strategy for feature selection. In: 9th Internat. Conf. on Pattern Recognition. IEEE, pp. 706–708.
- Yang, J., Honavar, V., 1998. Feature subset selection using a genetic algorithm. *IEEE Intell. Systems* 13, 44–49.
- Yu, B., Yuan, B., 1993. A more efficient branch and bound algorithm for feature selection. *Pattern Recognition* 26 (6), 883–889.