

Kernel-Based Feature Extraction with a Speech Technology Application

András Kocsor and László Tóth

Abstract—Kernel-based nonlinear feature extraction and classification algorithms are a popular new research direction in machine learning. This paper examines their applicability to the classification of phonemes in a phonological awareness drilling software package. We first give a concise overview of the nonlinear feature extraction methods such as kernel principal component analysis (KPCA), kernel independent component analysis (KICA), kernel linear discriminant analysis (KLDA) and kernel springy discriminant analysis (KSDA). The overview deals with all the methods in a unified framework, regardless of whether they are unsupervised or supervised. The effect of the transformations on a subsequent classification is tested in combination with learning algorithms such as Gaussian mixture modeling (GMM), artificial neural nets (ANN), projection pursuit learning (PPL), decision tree-based classification (C4.5) and support vector machines (SVM). We found in most cases that the transformations have a beneficial effect on the classification performance. Furthermore, the nonlinear supervised algorithms yielded the best results.

Index Terms—kernel feature spaces, kernel-based feature extraction, principal component analysis, independent component analysis, discriminant analysis, kernel-based methods

I. INTRODUCTION

AUTOMATIC speech recognition is dealt with by quite traditional statistical modelling techniques such as Gaussian mixture modelling or artificial neural nets. In the last couple of years, however, the theory of machine learning has developed a wide variety of novel learning and classification algorithms. In particular, the so-called kernel-based methods have recently become a flourishing new research direction. Kernel-based classification and regression techniques, including the well-known Support Vector Machines (SVM) found their way into speech recognition relatively slowly. This is probably because their application to such large-scale tasks as speech recognition required addressing both theoretical and practical problems. Recently, however, more and more authors are turning their attention to the application of SVM in speech recognition (see, e.g. [10], [17], [46], [50]).

Besides using kernel-based classifiers, an alternative option is to use kernel-based technologies only to transform the feature space, and leave the job of classification to more traditional methods [39]. The goal of this paper is to study the applicability of some of these methods to phoneme classification, making use of kernel-based feature extraction methods applied prior to learning in order to improve classification rates. In essence, this article deals with kernel principal component analysis (KPCA) [47], kernel independent component

analysis (KICA) [35], [3], kernel linear discriminant analysis (KLDA) [40], [5], [34] and kernel springy discriminant analysis (KSDA) [36] techniques. Their effect on classification performance is then tested in combination with classifiers such as Gaussian mixture modeling (GMM), artificial neural networks (ANN), projection pursuit learning (PPL) decision tree-based classification (C4.5) and support vector machines (SVM). The algorithms are applied to two recognition tasks. One of them is real-time phoneme recognition which, when combined with a real-time visualization of the results, forms the basis of the “SpeechMaster” phonological awareness drilling software developed by our team. The other test set was the well-known TIMIT phone classification task.

The structure of the paper is as follows. First, we provide a concise overview of the nonlinear feature extraction methods. The overview is written so that it deals with all the methods in a unified framework, regardless of whether they are unsupervised or supervised. Furthermore, the traditional linear counterparts of the methods can be obtained as special cases of our approach. Afterwards, we present the goals of the “SpeechMaster” software along with the phoneme classification problem that arises. Besides the special vowel recognition task the “SpeechMaster” software is built on we also present test results on the TIMIT database. In both cases, we first briefly describe the acoustic features that were applied in the experiments and also list the learning methods used. Then, in the final part of the paper we present the results of the experiments and discuss them from several aspects, focusing on the advantages and drawbacks of each nonlinear feature extraction method.

II. KERNEL-BASED FEATURE EXTRACTION

Classification algorithms require that the objects to be classified are represented as points in a multidimensional feature space. However, before executing a learning algorithm, additional vector space transformations may be applied on the initial features. The reason for doing this is twofold. Firstly they can improve classification performance and, secondly, they can reduce the dimensionality of the data. In the literature sometimes both the choice of the initial features and their transformation are dealt with under the name “feature extraction”. To avoid any misunderstanding, in this section it will cover only the latter, that is the transformation of the initial feature set into another one, which is hoped will yield a more efficient or, at least, faster classification.

The approach of feature extraction could be either linear or nonlinear, but there is a technique (which is most topical nowadays) that is, in some sense, breaking down the barrier

A. Kocsor and L. Tóth are with the Research Group on Artificial Intelligence of the Hungarian Academy of Sciences and University of Szeged, H-6720 Szeged, Aradi vértanúk tere 1., Hungary, E-mail: kocsor@inf.u-szeged.hu, tothl@inf.u-szeged.hu

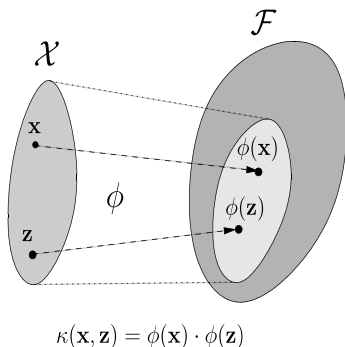


Fig. 1. The "kernel-idea". The dot product in the kernel feature space \mathcal{F} is defined implicitly.

between the two types. The key idea behind the kernel technique was originally presented in [1] and was again applied in connection with the general purpose Support Vector Machine [8], [11], [49], [52], [53], which was later followed by other kernel-based methods [3], [5], [34], [35], [36], [40], [45], [47], [48]. In the following we summarize four nonlinear feature extraction methods that may be derived using the kernel-based nonlinearizations of the linear algorithms PCA, ICA, LDA and SDA. We do not present the linear techniques separately because the non-linear descriptions will be formalized in such a way that, with a proper parametrization, they lead back to the traditional linear methods. All methods will be dealt with in a unified, concise form. We also represent the effects of these transformations on artificial data sets via figures. In addition, we always give references to sources where a detailed description on the feature extraction technique in question can be found. Our main aim is to help the reader gain a unified view of the methods and get some ideas about their usage. First, we provide a set of definitions. Then we discuss the kernel idea, followed by an explanation of each method one after the other. The section rounds off with some remarks about how the number of required calculations can be reduced by decimating the sample set.

A. Introduction

Without loss of generality we shall assume that as a realization of multivariate random variables, there are n -dimensional real attribute vectors in a compact set \mathcal{X} over \mathbb{R}^n describing objects in a certain domain, and that we have a finite $n \times k$ sample matrix $X = [\mathbf{x}_1, \dots, \mathbf{x}_k]$ containing k random observations. Let us assume as well that we have r classes and an indicator function

$$\mathcal{L} : \{1, \dots, k\} \rightarrow \{1, \dots, r\}, \quad (1)$$

where $\mathcal{L}(i)$ gives the class label of the sample \mathbf{x}_i . Let k_j further denote the number of vectors associated with label j in the sample data¹.

Now we continue with the definition of the kernel-based feature extraction and then outline the kernel idea.

¹The two types of feature extraction methods (supervised or unsupervised) can be distinguished by whether they utilize an indicator function or not during the computation of the transformation parameters.

The goal of feature extraction is to find a mapping $h : \mathcal{X} \rightarrow \mathcal{Y}$ which leads to a new set of features that are optimal according to a given criterion. In the case of kernel-based feature extraction the mapping is nonlinear and has the following form:

$$\mathbf{z} \rightarrow A\kappa(X, \mathbf{z}), \quad \mathbf{z} \in \mathcal{X} \quad (2)$$

where A is a constant, real $m \times k$ matrix, the function $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is continuous, symmetric and positive definite (which is called a Mercer kernel), X is the sample matrix and $\kappa(X, \mathbf{z})$ is a short-hand notation for the vector $[\kappa(\mathbf{x}_1, \mathbf{z}), \dots, \kappa(\mathbf{x}_k, \mathbf{z})]^\top$. As can be seen in Eq. (2), linear combinations of the base functions $\kappa(\mathbf{x}_1, \mathbf{z}), \dots, \kappa(\mathbf{x}_k, \mathbf{z})$ give the components of the new feature vector $h(\mathbf{z})$. The criterion (and hence the calculation that leads to matrix A in Eq. (2)) is different for each method and the result depends on the data set X .

Based on Mercer's theorem [53], if κ is a Mercer kernel then a dot product space \mathcal{F} necessarily exists with a mapping $\phi : \mathcal{X} \rightarrow \mathcal{F}$ (see Fig. 1) such that

$$\forall \mathbf{x}, \mathbf{z} \in \mathcal{X} \quad \kappa(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z}). \quad (3)$$

Usually \mathcal{F} is called the kernel feature space and ϕ is the feature map. At this point we have two immediate consequences. When ϕ is the identity, the function $\kappa(\mathbf{x}, \mathbf{z}) = \mathbf{x} \cdot \mathbf{z}$ (the simple dot product over the space \mathcal{X}) is symmetric, continuous and positive definite, so it constitutes a proper Mercer kernel. Going the other way, when applying a general Mercer kernel we can assume a space \mathcal{F} over which we perform dot product calculations. This space and dot product calculations over it are defined only implicitly via the kernel function itself. The space \mathcal{F} and map ϕ may not be explicitly known. We need only define the kernel function, which then ensures an implicit evaluation. The construction of Mercer kernels, when such a mapping ϕ exists, is a non-trivial problem, but there are some possible candidates available (cf. [13], [18], [32]). From the functions available, the three most popular are the *polynomial kernel* κ_1 , the *Gaussian RBF kernel* κ_2 and the *rational quadratic kernel* κ_3 :

$$\begin{aligned} \kappa_1(\mathbf{x}, \mathbf{z}) &= (\mathbf{x} \cdot \mathbf{z} + 1)^d, & d &\in \mathbb{N}, \\ \kappa_2(\mathbf{x}, \mathbf{z}) &= \exp(-\|\mathbf{x} - \mathbf{z}\|^2/\gamma), & \gamma &\in \mathbb{R}_+, \\ \kappa_3(\mathbf{x}, \mathbf{z}) &= 1 - \|\mathbf{x} - \mathbf{z}\|^2/(\|\mathbf{x} - \mathbf{z}\|^2 + 1), & \sigma &\in \mathbb{R}_+. \end{aligned} \quad (4)$$

Let us suppose that we have chosen a specific kernel function along with a proper feature map ϕ and a kernel feature space \mathcal{F} . Then the nonlinear mapping of Eq. (2) can be written as:

$$\begin{aligned} A\kappa(X, \mathbf{z}) &= \begin{pmatrix} a_{11}\kappa(\mathbf{x}_1, \mathbf{z}) + \dots + a_{1k}\kappa(\mathbf{x}_k, \mathbf{z}) \\ \vdots \\ a_{m1}\kappa(\mathbf{x}_1, \mathbf{z}) + \dots + a_{mk}\kappa(\mathbf{x}_k, \mathbf{z}) \end{pmatrix} \\ &= \begin{pmatrix} a_{11}\phi(\mathbf{x}_1)^\top + \dots + a_{1k}\phi(\mathbf{x}_k)^\top \\ \vdots \\ a_{m1}\phi(\mathbf{x}_1)^\top + \dots + a_{mk}\phi(\mathbf{x}_k)^\top \end{pmatrix} \phi(\mathbf{z}). \end{aligned}$$

In the following we will denote the latter matrix by V . Notice here that V is constant, and its rows contain linear

combinations of the image of the data vectors in \mathcal{F} . This means that the transformation is linear in the kernel feature space. But because the feature map itself is nonlinear we obtain a nonlinear transformation of the sample points of the initial feature space \mathcal{X} .

All the algorithms that we are going to present in the following are linear mappings in the kernel feature space, the row vectors of matrix V being obtained by optimizing a different objective function $\tau(\mathbf{v})$, say. What is common in each case is that we will look for \mathbf{v} directions with large values of τ . Intuitively, if larger values of τ indicate better directions and the row vectors of V need to be independent in certain ways, choosing stationary points that have the m largest function values is a reasonable strategy. Obtaining the above stationary points of a general objective function is a difficult global optimization problem. But if τ is defined by a Rayleigh quotient formulae:

$$\tau(\mathbf{v}) = \frac{\mathbf{v}^\top \mathcal{B}_1 \mathbf{v}}{\mathbf{v}^\top \mathcal{B}_2 \mathbf{v}}, \quad (5)$$

the solution is easy and fast when formulated as a generalized eigenvalue problem ($\mathcal{B}_1 \mathbf{v} = \lambda \mathcal{B}_2 \mathbf{v}$). Actually, this approach offers a unified view of the feature extraction methods discussed in this paper.

B. Kernel Principal Component Analysis

Principal component analysis [29] is a ubiquitous unsupervised technique for data analysis and dimension reduction. To explain how its nonlinear version works [47]², let us first choose a kernel function κ for which

$$\kappa(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z}), \quad \mathbf{x}, \mathbf{z} \in \mathcal{X},$$

holds for a mapping $\phi: \mathcal{X} \rightarrow \mathcal{F}$. It is well-known that PCA looks for those directions of \mathcal{X} in which the variance of the data is large. We will do exactly the same, but in the kernel feature space \mathcal{F} . For this we define the objective function τ as

$$\tau(\mathbf{v}) = \frac{\mathbf{v}^\top \mathcal{C} \mathbf{v}}{\mathbf{v}^\top \mathbf{v}}, \quad \mathbf{v} \in \mathcal{F} \setminus \{\mathbf{0}\}, \quad (6)$$

where \mathcal{C} is the covariance matrix of the image of the sample $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_k)$:

$$\mathcal{C} = E\{(\phi(\mathbf{x}) - E\{\phi(\mathbf{x})\})(\phi(\mathbf{x}) - E\{\phi(\mathbf{x})\})^\top\}. \quad (7)$$

Now we define the Kernel-PCA transformation based on the stationary points of Eq. (6), which are given as the eigenvectors of the symmetric positive semidefinite matrix \mathcal{C} . However, since this matrix is of the form

$$\mathcal{C} = \sum_{i,j}^k c_{ij} \phi(\mathbf{x}_i) \phi(\mathbf{x}_j)^\top, \quad (8)$$

we can suppose the following equation holds during the analysis of the stationary points:

$$\mathbf{v} = \sum_{i=1}^k \alpha_i \phi(\mathbf{x}_i). \quad (9)$$

²The derivation presented here differs slightly from the one originally proposed by Schölkopf. But the result of the derivation is equivalent to the original.

We can arrive at this assumption in many ways, e.g. we can decompose an arbitrary vector \mathbf{v} into vectors $\mathbf{v}_1 + \mathbf{v}_2$, where \mathbf{v}_1 gives that component of \mathbf{v} which falls in $\text{span}(\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_k))$, while \mathbf{v}_2 gives the component perpendicular to it. Since $C\mathbf{v} = C(\mathbf{v}_1 + \mathbf{v}_2) = C\mathbf{v}_1$ and for the stationary points the eigenvalue-eigenvector equality $C\mathbf{v} = \lambda\mathbf{v}$ is satisfied, we find that the condition defined in Eq. (9) (i.e. $\mathbf{v}_2^\top \mathbf{v}_2 = 0$) does not restrict generality.

Based on the above assumption the variational parameters of τ can be the vector α instead of \mathbf{v} :

$$\tau(\alpha) = \frac{\left(\sum_{i=1}^k \alpha_i \phi(\mathbf{x}_i)^\top\right) \mathcal{C} \left(\sum_{j=1}^k \alpha_j \phi(\mathbf{x}_j)\right)}{\left(\sum_{i=1}^k \alpha_i \phi(\mathbf{x}_i)^\top\right) \left(\sum_{j=1}^k \alpha_j \phi(\mathbf{x}_j)\right)}. \quad (10)$$

It is easy to see that

$$\tau(\alpha) = \frac{\alpha^\top \frac{1}{k} K (I - \bar{1}) K \alpha}{\alpha^\top K \alpha}, \quad (11)$$

where $[K]_{ij} = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = \kappa(\mathbf{x}_i, \mathbf{x}_j)$ is a Gram matrix, I is the unit matrix and $[\bar{1}]_{ij} = 1/k$.

After differentiating (11) with respect to α we find that the stationary points are the solution vectors of the general eigenvalue problem

$$\frac{1}{k} K (I - \bar{1}) K \alpha = \lambda K \alpha, \quad (12)$$

which is equivalent to the problem

$$\frac{1}{k} (I - \bar{1}) K \alpha = \lambda \alpha. \quad (13)$$

Although the matrix $(I - \bar{1})K$ is not symmetric, its eigenvalues are real and non-negative, and those eigenvectors that have positive eigenvalues are orthogonal. In fact the best approach is to solve the following symmetric eigenproblem, where the positive eigenvalues and the corresponding eigenvectors are the same as those obtained from Eq. (13)

$$\frac{1}{k} (I - \bar{1}) K (I - \bar{1}) \alpha = \lambda \alpha. \quad (14)$$

Now let the m positive dominant eigenvalues of $\frac{1}{k} (I - \bar{1}) K (I - \bar{1})$ be denoted by $\lambda_1 \geq \dots \geq \lambda_m$ and the corresponding eigenvectors be $\alpha_1, \dots, \alpha_m$. Then the matrix A of the transformation we need (cf. Eq. (2)) can be calculated like so:

$$A = [\alpha_1, \dots, \alpha_m]^\top. \quad (15)$$

The effect of KPCA is demonstrated in Fig. 2. The data set of Fig. (A) was transformed using linear PCA, that is KPCA was performed using the kernel $\kappa(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$. The result is shown in Fig. (B). Evidently, the algorithm found the direction with the largest variance and chose it as the x -axis of the transformed data. This effect is also justified by the shape of distribution curves shown below the images. In a second experiment the data set of Fig. (D) was transformed, but in this case using the rational quadratic kernel, which leads to a nonlinear transformation. The result is shown in Fig. (E). Examining the distribution of the points along the x -axis, one can see that the variance of the data has significantly increased owing to the nonlinearity of the method employed.

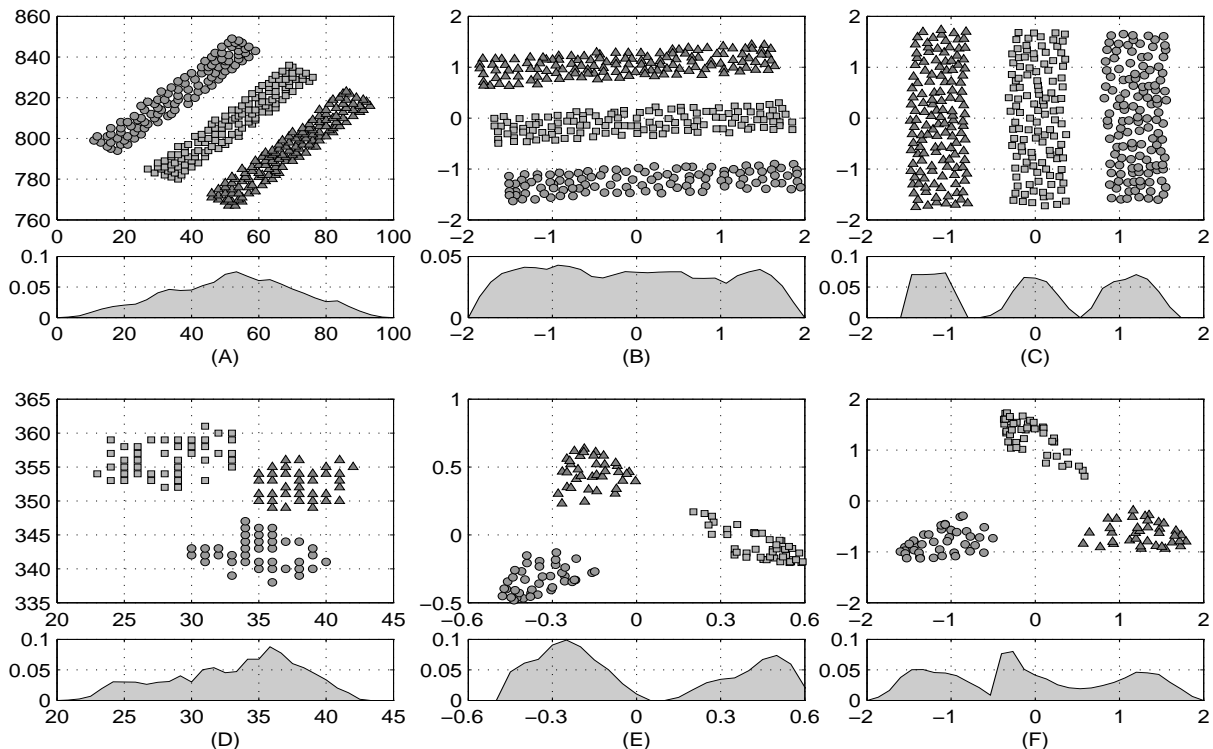


Fig. 2. This illustrates the typical behavior of KPCA and KICA. Figures (A) and (D) show some artificial data sets before the transformation. Figures (B) and (E) show the resulting distribution after linear and nonlinear KPCA, respectively. Figs. (C) and (F) depict the results of a linear and nonlinear KICA. The distribution of the data points along the x -axis is shown below each figure.

C. Kernel Independent Component Analysis

Independent Component Analysis [12], [15], [26], [27], [28] is a general-purpose statistical method that originally arose from the study of blind source separation (BSS). Another application of ICA is unsupervised feature extraction, where the aim is to linearly transform the input data into uncorrelated components, along which the distribution of the sample set is the least Gaussian. The reason for this is that along these directions the data is supposedly easier to classify. This is in concordance with the most common speech modeling technique, that is fitting mixtures of Gaussians on each class. Obviously, this assumes that the class distributions can be well approximated by Gaussian mixtures. ICA extends this by assuming that the distribution when all classes are fused, on the contrary, is not Gaussian, so, using non-Gaussianity as a heuristic for unsupervised feature extraction will prefer those directions that separate the classes.

For optimal selection of the independent directions several objective functions were defined using approximately equivalent approaches. The goal of the ICA algorithm itself is to find the optimum of these objective functions. There are many iterative methods for performing Independent Component Analysis. Some of these do require preprocessing, i.e. centering and whitening, while others do not. In general, experience shows that all these algorithms should converge faster on centered and whitened data, even with those which do not really require it.

Let us first examine how the centering and whitening preprocessing steps can be performed in the kernel feature

space. To this end, let the inner product be implicitly defined by the kernel function κ in \mathcal{F} with associated transformation ϕ .

Centering in \mathcal{F} . We shift the data $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_k)$ with its mean $E\{\phi(\mathbf{x})\}$, to obtain data

$$\begin{aligned} \phi'(\mathbf{x}_1) &= \phi(\mathbf{x}_1) - E\{\phi(\mathbf{x})\} \\ &\vdots \\ \phi'(\mathbf{x}_k) &= \phi(\mathbf{x}_k) - E\{\phi(\mathbf{x})\} \end{aligned} \quad (16)$$

with a mean of $\mathbf{0}$.

Whitening in \mathcal{F} . The goal of this step is to transform the centered samples $\phi'(\mathbf{x}_1), \dots, \phi'(\mathbf{x}_k)$ via an orthogonal transformation Q into vectors $\hat{\phi}(\mathbf{x}_1) = Q\phi'(\mathbf{x}_1), \dots, \hat{\phi}(\mathbf{x}_k) = Q\phi'(\mathbf{x}_k)$ where the covariance matrix $\hat{C} = E\{\hat{\phi}(\mathbf{x})\hat{\phi}(\mathbf{x})^\top\}$ is the unit matrix. Since standard principal component analysis [29] – just like its kernel-based counterpart – transforms the covariance matrix into a diagonal form, where the diagonal elements are the eigenvalues of the data covariance matrix $E\{\phi'(\mathbf{x})\phi'(\mathbf{x})^\top\}$, it only remains to transform each diagonal element to 1. Based on this observation, the required whitening transformation is obtained by slightly modifying the formulas presented in the section on KPCA. Now if we assume that the eigenpairs of $E\{\phi'(\mathbf{x})\phi'(\mathbf{x})^\top\}$ are $(\alpha_1, \lambda_1), \dots, (\alpha_k, \lambda_k)$ and $\lambda_1 \geq \dots \geq \lambda_k$, the transformation matrix Q will take the form $[\lambda_1^{-1/2}\alpha_1, \dots, \lambda_m^{-1/2}\alpha_m]^\top$. If m is less than n a dimensionality reduction is employed.

After the nonlinear preprocessing we can apply one of the many linear ICA algorithms. We present here the FastICA algorithm of Hyvarinen, for which centralization and whitening is a prerequisite.

For the sake of simplicity here we will denote the preprocessed data samples by $\mathbf{z}_1, \dots, \mathbf{z}_k$. In this new linear space we are going to search for directions along which the distribution of the data is the least Gaussian. To measure this criterion we introduce the following objective function:

$$J_G(\eta) \approx (E\{G(\eta)\} - E\{G(\nu)\})^2 \quad (17)$$

where η is a variable with zero mean and unit variance, $G: \mathbb{R} \rightarrow \mathbb{R}$ is an appropriate non-quadratic function, E again denotes the expectation value and ν is a standardized Gaussian variable. The following three choices of G are conventionally used:

$$\begin{aligned} G_1(\eta) &= \eta^4, \\ G_2(\eta) &= \log(\cosh(\eta)), \\ G_3(\eta) &= -\exp(-\eta^2/2). \end{aligned} \quad (18)$$

It should be mentioned here that in Eq. (17) the expectation value of $G(\nu)$ is a constant, its value only depending on the selected function (e.g. $E\{G_1(\nu)\} = 3$). The variable η has a leptokurtic distribution (a distribution with a high peak) if $E\{G_1(\eta)\} > 3$, it is a mesokurtic variable if $E\{G_1(\eta)\} = 3$, while it has platykurtic distribution (i.e. it is a flat-topped curve) when $E\{G_1(\eta)\} < 3$. For leptokurtic independent components the optimal contrast function is one which grows slower than quadratically, while the optimal for platykurtic components grows faster (cf. [28]). In Hyvärinen's FastICA algorithm for selecting a new direction \mathbf{v} the following τ objective function is used:

$$\tau_G(\mathbf{v}) = (E\{G(\mathbf{v}^\top \mathbf{z})\} - E\{G(\nu)\})^2, \quad (19)$$

which may be obtained by replacing η in Eq. (17) with $\mathbf{v}^\top \mathbf{z}$, the dot product of the direction \mathbf{v} and sample \mathbf{z} . FastICA is an approximate Newton iteration procedure for the local optimization of the function $\tau_G(\mathbf{v})$.

Before discussing the optimization problem, let us first examine the properties of the preprocessed data $\mathbf{z}_1, \dots, \mathbf{z}_k$.

- For every normalized \mathbf{v} vector the mean of $\mathbf{v}^\top \mathbf{z}_1, \dots, \mathbf{v}^\top \mathbf{z}_k$ is set to zero, and its variance is set to one. Actually we need this since Eq. (17) requires that η should have a zero mean and variance of one hence, with the substitution $\eta = \mathbf{v}^\top \mathbf{z}$, the projected data $\mathbf{v}^\top \mathbf{z}$ must also have this property.
- For any matrix W the covariance matrix C_W of the transformed preprocessed points $W\mathbf{z}_1, \dots, W\mathbf{z}_k$ will remain a unit matrix if and only if W is orthogonal, since

$$\begin{aligned} \hat{C}_W &= E\{W\mathbf{z}(W\mathbf{z})^\top\} = WE\{\mathbf{z}\mathbf{z}^\top\}W^\top \\ &= WIW^\top = WW^\top. \end{aligned} \quad (20)$$

After preprocessing, FastICA looks for a new orthogonal base W for the preprocessed data, where the values of the non-Gaussianity measure τ_G for the base vectors are large. Note that since the data remains whitened after an orthogonal transformation, ICA can be considered an extension of PCA.

Now we briefly outline how the FastICA algorithm works (c.f. [15], [27]). The input for this algorithm is the preprocessed sample $Z = [\mathbf{z}_1, \dots, \mathbf{z}_k]$ and the nonlinear function G , while the output is the transformation matrix W . The first and second order derivatives of G are denoted by G' and G'' .

```

procedure FastICA(Z, G);
% initialization
let  $W_0$  be a random  $m \times m$  matrix;
 $W_0 = (W_0 W_0^\top)^{-1/2} W_0$ ;
 $i = 0$ ;
% approximate Newton iteration
While  $W$  has not converged;
  for  $j = 1$  to  $m$ 
    let  $\mathbf{s}_j$  be the  $j$ th row vector of  $W_i$ ;
     $\mathbf{w}_j = E\{\mathbf{z}G'(\mathbf{s}_j \cdot \mathbf{z})\} - E\{G''(\mathbf{s}_j \cdot \mathbf{z})\}\mathbf{s}_j$ ;
  end;
   $i = i + 1$ ;
   $W_i = [\mathbf{w}_1, \dots, \mathbf{w}_m]^\top$ ;
   $W_i = (W_i W_i^\top)^{-1/2} W_i$ ;
do
End procedure

```

In the pseudo-code $(W_i W_i^\top)^{-1/2} W_i$ means a symmetric decorrelation, where $(W_i W_i^\top)^{-1/2}$ can be readily obtained from its eigenvalue decomposition. If $W_i W_i^\top = V D V^\top$, then $(W_i W_i^\top)^{-1/2}$ is equal to $V D^{-1/2} V^\top$. Finally, the expected values required by the algorithm are calculated as the empirical means of the preprocessed input samples in Z .

We should remark that in the discussion above we nonlinearized only centering and whitening, not the consecutive iterative FastICA algorithm. This would also be possible, as in τ_G the dot product $\mathbf{v}^\top \mathbf{z}$ could be nonlinearized with the kernel method. But this would go outside our unified discussion based on the Rayleigh quotient. Practically speaking, the Kernel FastICA method = Kernel-Centering + Kernel-Whitening + iterative process of the original FastICA. The transformation matrix A (cf. Eq. (2)) of KICA is WQ , where Q represents centering and whitening, while W corresponds to the orthogonal matrix produced by FastICA. Despite the fact that the second, optimization phase for finding W is not based on the Rayleigh quotient approach, we feel that KICA as a unique extension of KPCA can be the part of this review. More details on the family of the KICA methods can be found in [34], [3].

To demonstrate the behavior of KICA we return to the artificial data set in Fig. 2. We once again transformed the data sets (A) and (D) but now with KICA. Fig. (C) shows the result when using a linear kernel, while Fig. (F) shows the effect of a rational quadratic kernel. When compared with KPCA, it can be readily seen that while KPCA looks for directions with a large variance, KICA prefers those directions with the least possible Gaussian distribution.

D. Kernel Linear Discriminant Analysis

Linear Discriminant Analysis (LDA) is a traditional supervised feature extraction method [16] which has proved to be one of the most successful preprocessing techniques for classification. It has long been used in speech recognition as well [22], [4], [51]. The goal of LDA is to find a new (not necessarily orthogonal) basis for the data which provides the optimal separation between classes. To present the steps of KLDA we virtually follow the discussion of its linear counterpart, but in this case everything is meant to happen implicitly, in the kernel feature space \mathcal{F} .

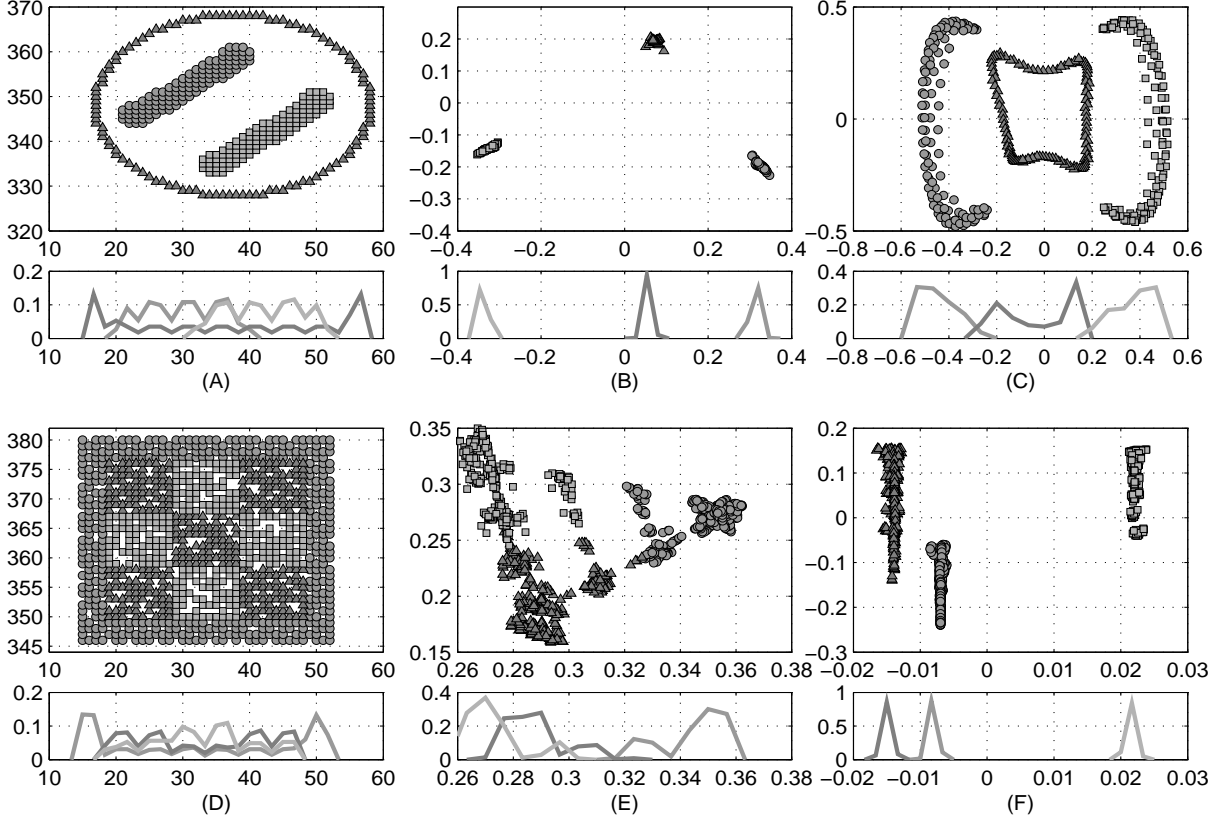


Fig. 3. This illustrates the effect of the supervised algorithms KLDA and KSDA. Figures (A) and (D) depict artificial data sets. Figures (B) and (E) show the resulting data sets after applying KLDA on (A) and (D), respectively. Figures (C) and (F) represent the KSDA-transformed versions of (A) and (D). The distributions of the classes along the x -axis is also shown below the figures. In every case the transformation applied was nonlinear.

Let us again suppose that a kernel function κ has been chosen along with a feature map ϕ and a kernel feature space \mathcal{F} . In order to define the transformation matrix A of KLDA we first define the objective function $\tau : \mathcal{F} \rightarrow \mathbb{R}$ which depends not only on the sample data X , but also on the indicator function \mathcal{L} owing to the supervised nature of this method. Let us define

$$\tau(\mathbf{v}) = \frac{\mathbf{v}^\top \mathcal{B} \mathbf{v}}{\mathbf{v}^\top \mathcal{W} \mathbf{v}}, \quad \mathbf{v} \in \mathcal{F} \setminus \{\mathbf{0}\}, \quad (21)$$

where \mathcal{B} is the *Between-class Scatter Matrix*, while \mathcal{W} is the *Within-class Scatter Matrix*. Here the *Between-class Scatter Matrix* \mathcal{B} shows the scatter of the class mean vectors $\boldsymbol{\mu}_j$ around the overall mean vector $\boldsymbol{\mu}$:

$$\begin{aligned} \mathcal{B} &= \sum_{j=1}^r \frac{k_j}{k} (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^\top \\ \boldsymbol{\mu} &= \frac{1}{k} \sum_{i=1}^k \phi(\mathbf{x}_i) \\ \boldsymbol{\mu}_j &= \frac{1}{k_j} \sum_{\mathcal{L}(i)=j} \phi(\mathbf{x}_i) \end{aligned} \quad (22)$$

The *Within-class Scatter Matrix* \mathcal{W} represents the weighted average scatter of the covariance matrices \mathcal{C}_j of the sample vectors with the class label j :

$$\begin{aligned} \mathcal{W} &= \sum_{j=1}^r \frac{k_j}{k} \mathcal{C}_j \\ \mathcal{C}_j &= \frac{1}{k_j} \sum_{\mathcal{L}(i)=j} (\phi(\mathbf{x}_i) - \boldsymbol{\mu}_j)(\phi(\mathbf{x}_i) - \boldsymbol{\mu}_j)^\top \end{aligned} \quad (23)$$

$\tau(\mathbf{v})$ is large when its nominator is large and its denominator is small or, equivalently, when in the kernel feature space \mathcal{F} the within-class averages of the sample projected onto \mathbf{v} are

far from each other and the variance of the classes is small. The larger the value of $\tau(\mathbf{v})$ the farther the classes will be spaced and the smaller their spreads will be.

We may also suppose without loss of generality here that $\mathbf{v} = \sum_{i=1}^k \alpha_i \phi(\mathbf{x}_i)$ holds during the search for the stationary points of Eq. (21). With this assumption, after some algebraic rearrangement we obtain the formula

$$\tau(\mathbf{v}) = \frac{\mathbf{v}^\top \mathcal{B} \mathbf{v}}{\mathbf{v}^\top \mathcal{W} \mathbf{v}} = \frac{\boldsymbol{\alpha}^\top K (R - \bar{1}) K \boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top K (I - R) K \boldsymbol{\alpha}} = \tau(\boldsymbol{\alpha}), \quad (24)$$

where K is the kernel matrix, $[\hat{1}]_{ij} = 1/k$ and

$$[R]_{ij} = \begin{cases} \frac{1}{k_t} & \text{if } t = \mathcal{L}(i) = \mathcal{L}(j) \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

This means that Eq. (21) can be expressed as dot products of $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_k)$ and that the stationary points of this equation can be computed using the real eigenvectors of $(K\Theta_2 K)^{-1} K\Theta_1 K$. Since in general $K\Theta_2 K$ is a positive semidefinite matrix, it can be forced to be invertible if we add a small positive constant ϵ to its diagonal, that is we work with $K\Theta_2 K + \epsilon I$ instead of $K\Theta_2 K$. This matrix is guaranteed to be positive definite and hence should always be invertible. This small act of cheating can have only a negligible effect on the stationary points of eq. (21). If we further assume that the real eigenvectors with the largest $m (< r)$ real eigenvalues of

$(K\Theta_2K)^{-1}K\Theta_1K$ are $\alpha_1, \dots, \alpha_m$, then the transformation matrix A (cf. Eq. (2)) will be $[\alpha_1, \dots, \alpha_m]^\top$.

The behavior of KLDA is illustrated in Fig. 3, in the two examples of (A) and (D). In both cases the application of the exponential kernel resulted in a nonlinear transformation that minimized the variance of the classes, while giving the best spatial class separation at the same time. The results are shown in Figs. (B) and (E), respectively. Noting the distribution of the classes along the x -axis, one can see that their separability has increased.

E. Kernel Springy Discriminant Analysis

As was shown in the previous section, the KLDA criterion leads to a non-symmetric matrix, the eigenvectors of which are not necessarily orthogonal. Furthermore, we had to apply the shifting of the eigenspectrum to avoid numerical complications during inversion. These issues give rise to the need for an objective function τ , which results in a supervised transformation and yields similar results to KLDA, but is orthogonal and avoids the numerical problems mentioned.

Now let the dot product be implicitly defined (see Fig. 1) by the kernel function κ in the kernel feature space \mathcal{F} with associated transformation ϕ :

$$\kappa(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{y}). \quad (26)$$

The name Kernel Springy Discriminant Analysis stems from the utilization of a spring & antisppring model, which involves searching for directions with optimal potential energy using attractive and repulsive forces. In our case sample pairs in each class are connected by springs, while those of different classes are connected by antispings. New features can be easily extracted by taking the projection of a new point in those directions having a small spread in each class, while different classes are spaced out as much as possible. Let $\delta(\mathbf{v})$, the potential of the spring model along the direction \mathbf{v} in \mathcal{F} , be defined by

$$\sum_{i,j=1}^k ((\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))^\top \mathbf{v})^2 [\Theta]_{ij}, \quad (27)$$

where

$$[\Theta]_{ij} = \begin{cases} -1, & \text{if } \mathcal{L}(i) = \mathcal{L}(j) \\ 1, & \text{otherwise} \end{cases} \quad i, j = 1, \dots, k. \quad (28)$$

Naturally, the elements of matrix Θ can be initialized with values different from ± 1 as well. Each element of the matrix can be considered as a kind of spring quotient and each can be set to a different value for any pair of data points.

As before, we again suppose that the directions \mathbf{v} can be constructed as the linear combinations of the images of the data points in \mathcal{F} . That is,

$$\mathbf{v} = X_\phi \alpha \quad (29)$$

where $X_\phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_k)]$. To find the directions with large potentials, let the objective function τ be defined by

$$\tau(\alpha) = \frac{\delta(X_\phi \alpha)}{\alpha^\top \alpha}. \quad (30)$$

It is easy to prove that $\tau(\alpha)$ is equal to the following Rayleigh quotient formula

$$\tau(\alpha) = \frac{\alpha^\top X_\phi^\top \mathcal{B} X_\phi \alpha}{\alpha^\top \alpha}, \quad (31)$$

where

$$\mathcal{B} = \frac{1}{2} \sum_{i,j=1}^k (\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)) (\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j))^\top [\Theta]_{ij}. \quad (32)$$

Moreover, it is also straightforward to prove that Eq. (31) takes the following form:

$$\frac{\alpha^\top (K\tilde{\Theta}K - K\Theta K) \alpha}{\alpha^\top \alpha}, \quad (33)$$

where K is again the kernel matrix and $\tilde{\Theta}$ is a diagonal matrix with the sum of each row of Θ in the diagonal. After taking the derivative of Eq. (33) it is readily seen that the stationary points of $\tau(\alpha)$ can be obtained via an eigenanalysis of the following symmetric eigenproblem:

$$(K\tilde{\Theta}K - K\Theta K)\alpha = \lambda\alpha. \quad (34)$$

If we assume that the dominant m eigenvectors are $\alpha_1, \dots, \alpha_m$ then the transformation matrix A in Eq. (2) is defined by $[\alpha_1, \dots, \alpha_m]^\top$.

The effect of KSDA can again be visualized by transforming the data sets (A) and (D) of Fig. 3. While KLDA aims at minimizing the within-class variance and maximizing the between-class distance, KSDA does something similar, but based on within-class attractive and between-class repulsive forces. The results presented in Figs. (C) and (F) have a clearly separable class structure like those obtained using KLDA.

F. Reducing the Computational Cost

As we have already seen, all four methods lead to a (generalized) eigenproblem that involves finding the stationary points of the objective function $\tau(\mathbf{v})$, defined in the form of a Rayleigh quotient. During optimization, the vector \mathbf{v} consists of the linear combinations of the images of the data points X in the kernel feature space. Without doubt, if the amount of data points (k) is large, then the $k \times k$ sized matrices that are needed for constructing $\tau(\mathbf{v})$ – hence for solving the eigenproblem – can be so big that they pose serious computational and memory management problems.

Fortunately, in most practical problems good \mathbf{v} directions can be found even if we use only $r \ll k$ data points instead of k when constructing the linear combinations. Let us denote the indices of these r samples by $1 \leq i_1 < \dots < i_r \leq k$. It is easy to check that by just using these data items the formulas we obtain for the function $\tau(\alpha)$ can be expressed by the following:

$$\text{KPCA, KICA: } \frac{\alpha^\top \frac{1}{k} K_1^\top (I - \bar{1}) K_1 \alpha}{\alpha^\top K_2 \alpha}, \quad (35)$$

$$\text{KLDA: } \frac{\alpha^\top K_1^\top (R - \bar{1}) K_1 \alpha}{\alpha^\top K_1^\top (I - R) K_1 \alpha}, \quad (36)$$

$$\text{KSDA: } \frac{\alpha^\top (K_1^\top (\tilde{\Theta} - \Theta) K_1) \alpha}{\alpha^\top \alpha}, \quad (37)$$

where α is a vector of dimension r , K_1 is the matrix constructed from the columns i_1, \dots, i_r of the kernel matrix K , and K_2 is the minor matrix determined by the rows and columns of K with indices i_1, \dots, i_r . Based on these formulas, the eigenproblems to be solved are now reduced to a matrix of size $r \times r$. In practice, this matrix usually has no more than a couple of dozen or a couple of hundred rows and columns.

Of course, a key issue here is the strategy for choosing the r indices. Numerous selection strategies are possible from the random selection to the exhaustive search approach. In this paper we restrict our investigations to two different selection techniques. The first one is the simplest case when we chose samples randomly, why in the second we employed the kernel variant of the sequential forward floating selection (SFFS [43]) method with the LDA optimization criterion [37].

One more issue occurs that we need to discuss here. It is well-known that for the linear feature extraction methods PCA, ICA, LDA and SDA, the size of the problem is that of the original feature space. However, it depends on the number of the samples in the kernel counterparts. Despite these differences, if the kernel function is defined by the simple dot product $\kappa(\mathbf{x}, \mathbf{z}) = \mathbf{x}^\top \mathbf{z}$ and the feature map ϕ is realized by the identity ($\phi(\mathbf{x}) = \mathbf{x}$) then the kernel formulation of the methods (dual representation) are undoubtedly equivalent to the corresponding linear cases (primal representation). Obviously, as in practice the feature space is of lower dimension, it is worth using the linear methods when the simple dot product kernel is chosen. Now we show that the nonlinear formulae obtained this kernel function are readily traced back to the linear ones. Let us notice that in this case the kernel matrix K is equal to $X^\top X$, thus

$$\text{PCA, ICA: } \frac{\alpha^\top \frac{1}{k} X^\top X (I - \bar{1}) X^\top X \alpha}{\alpha^\top X^\top X \alpha} = \frac{\mathbf{v}^\top \frac{1}{k} X (I - \bar{1}) X^\top \mathbf{v}}{\mathbf{v}^\top \mathbf{v}}, \quad (38)$$

$$\text{LDA: } \frac{\alpha^\top X^\top X (R - \bar{1}) X^\top X \alpha}{\alpha^\top X^\top X (I - R) X^\top X \alpha} = \frac{\mathbf{v}^\top X (R - \bar{1}) X^\top \mathbf{v}}{\mathbf{v}^\top X (I - R) X^\top \mathbf{v}}, \quad (39)$$

$$\text{SDA: } \frac{\alpha^\top X^\top X (\tilde{\Theta} - \Theta) X^\top X \alpha}{\alpha^\top \alpha} = \frac{\mathbf{v}^\top X (\tilde{\Theta} - \Theta) X^\top \mathbf{v}}{\mathbf{v}^\top \mathbf{v}}, \quad (40)$$

where the vector $\mathbf{v} = X\alpha$, and matrices $X(I - \bar{1})X^\top$, $X(R - \bar{1})X^\top$, $X(I - R)X^\top$ and $X(\tilde{\Theta} - \Theta)X^\top$ are of the lower dimension.

III. EXPERIMENT NO.1.: CLASSIFICATION OF STEADY-STATE VOWELS

A. Application: A Phonological Awareness Teaching System

The ‘‘SpeechMaster’’ software developed by our team seeks to apply speech recognition technology to speech therapy and the teaching of reading. The role of speech recognition is to provide a visual phonetic feedback. In the first case it is intended to supplement the missing auditive feedback of the hearing impaired, while in the case of the latter it is to reinforce the correct association between the phoneme-grapheme pairs. With the aid of a computer children can practice without the need for the continuous presence of the teacher. This is very important because the therapy of the hearing impaired requires a long and tedious fixation phase. Furthermore, experience shows that most children prefer computer exercises to conventional drills.

Both applications require a real-time response from the system in the form of an easily comprehensible visual feedback. With the simplest display setting feedback is given by means of flickering letters, their identity and brightness being adjusted to the speech recognizer’s output. Figure 4 shows the user interface of ‘‘SpeechMaster’’, in the teaching reading and the speech therapy applications, respectively. As one can see, in the first case the flickering letter is positioned over a traditional picture for associating the word and word sound, while in the latter case it is combined with a web camera image which helps the impaired student learn the proper articulator positions.

B. Evaluation Domain

For training and testing purposes we recorded samples from 160 children aged between 6 and 8. The ratio of girls and boys was 50% - 50%. The speech signals were recorded and stored at a sampling rate of 22050 Hz in 16-bit quality. Each speaker uttered all the Hungarian vowels, one after the other, separated by a short pause. Since we decided not to discriminate their long and short versions, we only worked with 9 vowels altogether. The recordings were divided into a train and a test set in a ratio of 50% - 50%.

C. Acoustic Features

There are numerous methods for obtaining representative feature vectors from speech data [24], but their common property is that they are all extracted from 20-30 ms chunks or ‘‘frames’’ of the signal in 5-10 ms time steps. The simplest possible feature set consists of the so-called bark-scaled filterbank log-energies (FBLE). This means that the signal is decomposed with a special filterbank and the energies in these filters are used to parameterize speech on a frame-by-frame basis. In our tests the filters were approximated via Fourier analysis with a triangular weighting, as described in [24]. Altogether 24 filters were necessary to cover the frequency range from 0 to 11025 Hz. Although the resulting log-energy values are usually sent through a cosine transform to obtain the well-known mel-frequency cepstral coefficients, we abandoned it for two reasons: (I) the transforms we were going to apply have a similar decorrelating effect, and (II) we observed earlier that the learners we work with - apart from GMM - are not sensitive to feature correlation so, consequently, the cosine transform would bring no significant improvement [33]. Furthermore, as the data consisted of steady-state vowels, we found in a pilot test that adding the usual delta and delta-delta features could only marginally improve the results. So only the 24 filter bank log-energies formed this feature set, always extracted from the center frame of the vowels. Although it would be possible to stack several neighboring frames to form a larger feature set, because of the special steady-state nature of the vowel data used we saw no point in doing so.

The filterbank log-energies seem to be a proper feature set for a general speech recognition task as their spectro-temporal modulation is supposed to carry all the speech information [41]. But in the special task of classifying vowels pronounced in isolation it is only the gross spectral shape that



Fig. 4. Screenshots of the ‘SpeechMaster’ phonological awareness teaching system. (A-B) The teaching reading part and the speech therapy part, respectively.

carries the phonetic information. More precisely, it is known from phonetics that the spectral peaks (called formants) code the identity of vowels [41]. To estimate the formants, we implemented a simple algorithm that calculates the gravity centers and the variance of the mass in certain frequency bands [2]. The frequency bands are chosen so that they cover the possible place of the first, second and third formants. This resulted in 6 new features altogether.

A more sophisticated option for the analysis of the spectral shape would be to apply some kind of auditory model [21]. Unfortunately, most of these models are too slow for a real-time application. For this reason we experimented with the In-Synchrony-Bands-Spectrum of Ghitza [19], because it is computationally simple and attempts to model the dominance relations of the spectral components. The model analyses the signal using a filterbank that is approximated by weighting the output of an FFT - quite similar to the FBLE analysis. In this case, however, the output is not the total energy of the filter, but the frequency of the component that has the maximal energy, and so dominates the given frequency band. Obviously, the output resulting from this analysis contains no information about the energies in the filters, but only about their relative dominance. Hence we supposed that this feature set complements the FBLE features in a certain sense.

D. Learners

Describing the mathematical background of the learning algorithms applied is beyond the scope of this paper. Besides, we believe that they are familiar to those who are acquainted with pattern recognition. So in the following we specify only the parameters and the training algorithms used with each learner, respectively.

1) *Gaussian mixture modeling*: The most widely used method for modeling the class-conditional (continuous) distribution of the features is to approximate it by means of a weighted sum of Gaussians [14]. Traditionally the parameters are optimized according to the Maximum Likelihood (ML) criterion, using the expectation-maximization (EM) algorithm. It is well known, however, especially in the speech community

that maximum likelihood training is not optimal from a discrimination point of view, as it disregards the competing classes. Several alternatives have been proposed, such as Maximum Mutual Information (MMI) [42], [54] or Minimum Classification Error (MCE) criteria [30], [31]. Although these alternative training methods can significantly boost the classification performance, the increased computational requirements – especially when embedded in a hidden Markov model (HMM) – seems to be a deterrent to their widespread usage. Here we will utilize the EM algorithm with the following setup. As EM is an iterative technique, it requires a proper initialization of the parameters. To find a good starting parameter set we applied k -means clustering [16]. Since k -means clustering again only guaranteed finding a local optimum, we ran it 15 times with random parameters and used the one with the highest log-likelihood to initialize the EM algorithm. After experimenting, the best value for the number of mixtures k was found to be 3. In all cases the covariance matrices were forced to be diagonal.

2) *Artificial neural networks*: Since it was realized that, under proper conditions, ANNs can model the class posteriors [7], neural nets are becoming evermore popular in the speech recognition community. In the ANN experiments we used the most common feed-forward multilayer perceptron network with the backpropagation learning rule. The number of neurons in the hidden layer was set at 18 in each experiment (this value was chosen empirically, based on preliminary experiments). Training was stopped based on the cross-validation of 15% of the training data.

3) *Projection Pursuit Learning*: Projection pursuit learning is a relatively little-known modelling technique. It can be viewed as a neural net where the rigid sigmoid function is replaced by an interpolating polynomial. With this modification the representation power of the model is increased, so less units are necessary. Moreover, there is no need for additional hidden layers: one layer plus a second layer with linear combinations will suffice. During learning the model looks for directions in which the projection of the data points can be well approximated by its polynomials, thus the mean

square error will have the smallest value (hence the name ‘projection pursuit’). Our implementation follows the paper of [25]. In each experiment, a model with 8 projections and a 5th-order polynomial was applied.

4) *Support Vector Machines*: Support vector machines is a classifier algorithm that is based on the same kernel idea that we presented earlier. It first maps the data points into a high-dimensional feature space by applying some kernel function. Then, assuming that the data points have become easily separable in the kernel-space, it performs linear classifications to separate the classes. A linear hyperplane is chosen with a maximal margin. For further details on SVM the reader may peruse [53]. In all the experiments with SVM the radial basis kernel function was applied.

E. Experimental Setup

In the experiments 5 feature sets were constructed from the initial acoustic features, as described in Section III.B. *Set1* contained the 24 FBLE features. In *Set2* we combined *Set1* with the gravity center features, so *Set2* contained 30 measurements. *Set3* was composed of the 24 SBS features, while in *Set4* we combined the FBLE and SBS sets. Lastly, in *Set5* we added all the FBLE, SBS and gravity center features, thus obtaining a set of 54 values.

As regards the transformations, in every case we kept only the first 8 components. We performed this severe dimension reduction in order to show that, when combined with the transformations, the classifiers can yield the same scores in spite of the reduced feature set. To study the effects of nonlinearity, the linear version of each transformation was also used on each feature set. To obtain a sparse data representation for the kernel methods, we reduced the number of data points to 200 by applying the SFFS selection technique discussed earlier. Preliminary experiments showed that using more data would have no significant effect on the results.

In the classification experiments every transformation was combined with every classifier on every feature set. This resulted in the large table of Table I. In the header of the table PCA, ICA, LDA and SDA stand for the linear transformations (i.e. the kernel $\mathbf{x}^T \mathbf{z}$ was used), while KPCA, KICA, KLDA and KSDA stand for the nonlinear transformations (with an exponential kernel), respectively. The numbers shown are the recognition errors on the test data. The number in parenthesis denotes the number of features preserved after transformation. The best scores of each set are given in bold.

F. Results and Discussion

Upon inspecting the results the first thing one notices is that the SBS feature set (*Set3*) did about twice as badly as the other sets, no matter what transformation or classifier was tried. When combined with the FBLE features (*Set1*) both the gravity center and the SBS features brought some improvement, but this improvement is quite small and varies from method to method.

When focusing on the performance of the classifiers, ANN, PPL and SVM yielded very similar results. They, however, consistently outperformed GMM, which is the method most

commonly used in speech technology today. Firstly, this can be attributed to the fact that the functions that a GMM (with diagonal covariances) is able to represent are more restricted in shape than those of ANN or PPL. Secondly, it is a consequence of modeling the classes separately, rather than in the case of the other three classifiers that optimize a discriminative error function.

As regards the transformations, an important observation is that after the transformations the classification scores did not get worse compared to the classifications when no transformation was applied. This is so in spite of the dimension reduction, which shows that the features are highly redundant. Removing this redundancy by means of a transformation can make the classification more robust and, of course, faster.

Comparing the linear and the kernel-based algorithms, there is a slight preference towards the supervised transformations rather than the unsupervised ones. Similarly, the nonlinear transforms yielded somewhat better scores than the linear ones. The best transformation-classifier combination, however, varies from set to set. This warns us that no such broad claim can really be made about one transformation being superior to the others. This is always dependent on the feature set and the classifier. This is, of course, in accordance with the ‘‘no free lunch’’ theorem which claims that, for different learning tasks, different inductive bias can be beneficial [14].

Finally, we should make some general remarks. First of all, we must emphasize that both the transformations and the classifiers have quite a few adjustable parameters, and to examine all parameter combinations is practically impossible. Changing some of these parameters can sometimes have a significant effect on the classification scores. Keeping this (and the no free lunch theorem) in mind, our goal in this paper was to show that the nonlinear supervised transformations have the tendency to perform better (with any given classifier) than the linear and/or unsupervised methods. The results here seem to justify our hypothesis.

IV. EXPERIMENT NO.2.: TIMIT PHONE CLASSIFICATION

A. Evaluation Domain

In the vowel experiments the database, the number of features and the number of classes were all smaller than in a common speech recognition task. To assess the applicability of the algorithms to larger-scale problems, we also ran phone classification tests on the TIMIT database. The train and test sentences were chosen as usual, that is 3696 ‘sx’ and ‘si’ sentences formed the train set (142909 phone instances), and the complete test set (1344 ‘si’ and ‘sx’ sentences) were used for testing (51681 phone instances). The phone labels were fused into 39 classes, according to [38].

B. Acoustic Features

For the frame-based description of the signals we again used the bark-scaled filterbank log-energies. 22 filters were applied to cover the 0-8000 Hz frequency range of the TIMIT recordings.

Because the phonetic segments of the corpus are composed of a varying number of frames, an additional step was required

TABLE I

RECOGNITION ERRORS ON EACH FEATURE SET AS A FUNCTION OF THE TRANSFORMATION AND CLASSIFICATION APPLIED.

feature set	classifier	none (all)	PCA (8)	ICA (8)	LDA (8)	SDA (8)	KPCA (8)	KICA (8)	KLDA (8)	KSDA (8)
Set1 (24)	GMM	16.38 %	13.81 %	16.45 %	14.37 %	15.06 %	15.20 %	13.68 %	12.43 %	12.70 %
	ANN	10.34 %	9.86 %	9.93 %	10.97 %	9.58 %	9.86 %	9.58 %	8.05 %	7.98 %
	PPL	11.04 %	10.06 %	10.69 %	9.51 %	9.93 %	8.95 %	9.51 %	7.98 %	8.75 %
	SVM	9.93 %	10.00 %	8.95 %	8.05 %	8.05 %	8.88 %	8.26 %	6.73 %	7.22 %
Set2 (30)	GMM	13.33 %	11.38 %	13.33 %	12.84 %	13.33 %	13.47 %	12.36 %	10.27 %	11.31 %
	ANN	7.43 %	8.05 %	7.36 %	7.77 %	6.18 %	6.52 %	8.19 %	5.69 %	6.66 %
	PPL	9.37 %	8.59 %	6.54 %	6.11 %	6.45 %	6.59 %	6.45 %	4.93 %	6.66 %
	SVM	8.33 %	6.66 %	6.66 %	6.45 %	5.13 %	7.36 %	6.11 %	5.27 %	5.34 %
Set3 (24)	GMM	25.90 %	23.19 %	25.90 %	22.91 %	24.37 %	25.13 %	24.65 %	23.05 %	21.45 %
	ANN	20.00 %	18.88 %	19.58 %	21.45 %	20.00 %	21.04 %	18.54 %	18.26 %	17.84 %
	PPL	20.48 %	20.69 %	19.58 %	20.00 %	20.76 %	18.88 %	19.16 %	17.84 %	18.54 %
	SVM	19.65 %	20.69 %	18.88 %	17.36 %	19.58 %	19.79 %	18.33 %	16.52 %	16.45 %
Set4 (48)	GMM	13.95 %	12.01 %	15.90 %	13.81 %	14.16 %	15.34 %	12.08 %	10.00 %	9.93 %
	ANN	10.27 %	9.86 %	8.05 %	9.02 %	8.95 %	7.36 %	9.86 %	5.55 %	7.56 %
	PPL	10.48 %	8.95 %	9.37 %	8.95 %	9.44 %	7.36 %	9.09 %	6.18 %	7.98 %
	SVM	9.09 %	9.79 %	8.26 %	6.04 %	7.56 %	8.75 %	5.97 %	5.76 %	6.25 %
Set5 (54)	GMM	15.48 %	12.29 %	13.33 %	11.04 %	13.75 %	11.73 %	11.87 %	10.83 %	11.59 %
	ANN	8.68 %	7.01 %	6.45 %	10.00 %	7.56 %	9.09 %	6.59 %	7.15 %	4.93 %
	PPL	8.26 %	9.23 %	7.36 %	6.52 %	7.29 %	8.05 %	7.77 %	6.18 %	7.77 %
	SVM	9.37 %	8.54 %	5.76 %	4.65 %	5.62 %	6.11 %	5.76 %	6.18 %	4.23 %

to make them tractable for the transformations and learners, as these need all segments to be represented by the same number of features. For this we applied the very simple strategy of dividing each segment into three thirds, and averaging the filterbank energies over these subsegments (from a signal processing view this means a non-uniform smoothing and resampling). This method was popularized mainly in the SUMMIT system [20], but was also successfully applied by others as well [10]. To allow the learner to model the observation context at least to a certain level, additional average filterbank energies were calculated at the beginning and end of the segments. For this aim 50-50 ms intervals were considered at both sides.

Besides the resulting $5 \times 22 = 110$ energy-based features per segment, the length of the phone was also utilized. Furthermore, similar to the usual frame-based description strategies, we found that derivative-like features can be very useful – but in our case extracted only at the segment boundaries. These were calculated by RASTA-filtering the energy trajectories and then simply taking the frame-based differences at the boundaries. The role of RASTA filtering is to smooth the trajectories by removing those modulation frequency components that are perceptually not important [23]. In preliminary experiments we have found that it is unnecessary to calculate these delta-features in every bark-wide frequency channel. Rather, we have concluded that it is enough to extract them from fewer but wider frequency bands (this idea was in fact motivated by physiological results on the tuning curves of cochlear nucleus onset cells). Accordingly, only four 6-bark wide channels were used to calculate the delta features, altogether resulting in 8 of them (4 - 4 at each boundaries).

Finally, we have observed that smoothing over the segment thirds can sometimes remove important information, especially when working with long phone instances. To alleviate this, we extended our feature set with the variances of the energies calculated over the segments. These were again calculated

only from the 4 wide bands described above. Altogether, 123 segmental features were extracted from every phone instance. To justify the correctness of our representation we ran some preliminary classification tests, and the results were very close to those of others using a similar feature extraction technique [20], [10].

C. Learners

The TIMIT data set is much bigger than our vowel database. Consequently, we had no capacity to test every combination of the classifiers and learners, as we did in the case of the vowel data. Thus we decided to restrict ourselves to two classifiers only. ANN was chosen because of its consistently good performance and relatively small training time. The other classifier was selected based on the following rationale. The main aim of transforming the features space is to rearrange the data points so that they become more easily modelable by the subsequent learner. In accordance with this, the transforms must bring the most improvement when applied prior to a learner with a relatively small representation power. So, as the second classifier we chose C4.5. This is a well-known classifier in machine learning and, when trained on numerical data, it has a very restricted representation technique.

1) *Artificial neural networks*: In all the experiments the ANN had 38 inputs and 300 neurons in the hidden layer. Training was stopped based on cross-validation over 15% of the training data.

2) *C4.5*: C4.5 is a very well-known and widely used classifier in the machine learning community [44]. For those who prefer a statistical view, very similar learning schemes can be found under the name Classification and Regression Trees [9]. This method builds a tree-based representation from the data, and was originally invented with nominal features in mind. The algorithm was however extended for the case of numerical features. In this case the algorithm decomposes the feature space into rectangular blocks by means of axis-wise

hyperplanes. The hypercubes are iteratively decomposed into smaller and smaller ones, according to an entropy-based tree-building rule. This hashing of the feature space can be stopped by many possible criteria. Finally class labels are attached to each hyperbox, but posterior probability estimations are also easily attainable based on frequency counts. Obviously, the limited representation power of the model is caused both by the axis-wise restriction on the hyperplanes and the step-like look of the resulting probability estimations. In the experiments we used the original implementation of Quinlan. During tree building the minimum number of data points per leave was set to 24. The default parameters were used in every other respect.

D. Experimental Setup

Both the ANN and C4.5 classifiers were combined with each transformation. In the case of the kernel algorithms we always used the Gaussian RBF kernel (see Eq. 4). The number of features extracted by the transformations was always set to 38, that is the number of classes minus one. This value was chosen because LDA cannot return any more components (without tricks like splitting each class into subclasses), and to keep the results comparable we used the same number of features for the other transformations as well. As regards sparse data representation, because of the large size of the database we could not apply the SFFS technique (as its memory requirement is a quadratic function of the database size). So we decided to select the data points randomly, starting from 100 points, and iteratively adding further sets of 100 points. This was done in order to see how the number of points affected performance.

We were also interested in whether the choice of the contrast function of ICA influences its class separation abilities. To this end to learn more about this we performed tests with all three contrast functions listed in Eq. (18). Both linear ICA and Kernel-ICA (with an RBF kernel) were tried with all three contrast functions. The results showed that there were only small differences, but on the TIMIT data the contrast function G_3 seemed to behave the best. In the rest of the test we always worked with this contrast function.

E. Results and Discussion

The results of iteratively increasing the number of data points are plotted in Fig. 5. On every set Kernel-LDA was applied, with a subsequent ANN and C4.5 learning. The diagram shows how the classification error changes when the number of data is increased with a step size of 100. Clearly, the improvement is more dramatic for the C4.5 than for the ANN. In both cases there was no significant improvement beyond a sample size of 600. In the following experiments we always used this set of 600 points in the kernel-based tests.

The classification errors are summarized in Table II. Independent of the learner applied, we can say that the supervised algorithms performed better than the unsupervised ones, and that the kernel-based methods outperformed their linear counterparts. The differences are more significant in the case of the C4.5 learner than in the case of ANN. This is obviously

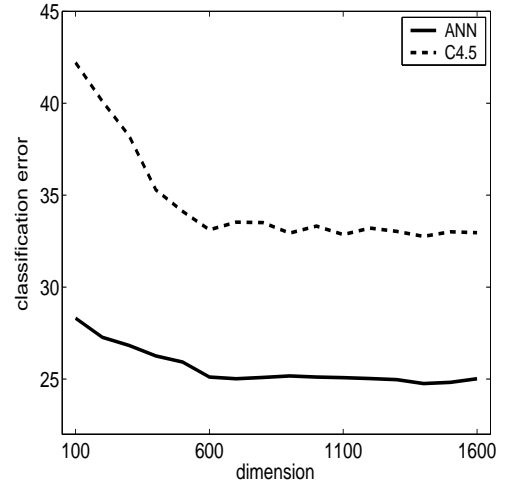


Fig. 5. Classification error as a function of the number of points kept in the sparse representation.

TABLE II
RECOGNITION ERRORS ON TIMIT

	ANN	C4.5		ANN	C4.5
none	28.43 %	49.86 %	none	28.43 %	49.86 %
PCA	28.22 %	47.23 %	KPCA	26.49 %	38.42 %
ICA	29.01 %	50.11 %	KICA	27.23 %	40.21 %
LDA	27.12 %	41.31 %	KLDA	25.11 %	33.12 %
SDA	27.93 %	43.12 %	KSDA	25.93 %	34.68 %

(A)

(B)

because of the flexibility of ANN representation, compared to the axis-wise rigid separation hyperplanes of C4.5.

V. CONCLUSIONS AND FUTURE WORK

The main purpose of this paper was to compare several classification and transformation methods applied to phoneme classification. The goal of applying a transformation can be dimension reduction, improvement of the classification scores, or increasing the robustness of the learning by removing the noisy and redundant features.

We found that nonlinear transformations in general lead to better classification than the nonlinear ones, and thus are a promising new direction for research. We also found that the supervised transformations are usually better than the unsupervised ones. We think that it would be worth looking for other supervised techniques that could be constructed in a similar way to the SDA or LDA technique. These transformations greatly improved our phonological awareness teaching system by offering a robust and reliable real-time phoneme classification. They also result in increased performance on the TIMIT data.

Finally, we should mention that finding the optimal parameters both for the transformations and the classifiers is quite a difficult problem. In particular, the parameters of the transformation and the subsequent learner are optimized separately at present. A combined optimization should probably produce better results, and there are already promising results in this direction in the literature [6]. Hence we plan to investigate parameter tuning and combined optimization.

REFERENCES

- [1] M. A. Aizerman, E. M. Braverman, L. I. Rozonoer, "Theoretical foundation of the potential function method in pattern recognition learning," *Automat. Remote Cont.*, Vol. 25, pp. 821-837, 1964.
- [2] D. Albesano, R. De Mori, R. Gemello, F. Mana, "A study on the effect of adding new dimensions to trajectories in the acoustic space," *Proc. of Eurospeech'99*, Budapest, pp. 1503-1506.
- [3] F. R. Bach, and M. I. Jordan, "Kernel Independent Component Analysis," *J. Machine Learning Res.* Vol. 3, pp. 1-48, 2002.
- [4] L.R. Bahl, P.V. deSouza P.S. Gopalakrishnan, D. Nahamoo and M. Picheny, "Robust Methods For Context Dependent Features And Models In A Continuous Speech Recognizer," *Proc. of ICASSP'94*, Adelaide, Australia, pp. 533-535., 1994.
- [5] G. Baudat and F. Anouar, "Generalized discriminant analysis using a kernel approach," *Neural Comput.*, Vol. 12, pp. 2385-2404, 2000.
- [6] A. Biem, S. Katagiri, E. McDermott and B.-H. Juang, "An application of discriminative feature extraction to filter-bank-based speech recognition." *IEEE Transactions on Speech and Audio Processing*, Vol. 9, no. 2, pp. 96-110., 2001.
- [7] C. M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press Inc., New York, 1996.
- [8] B. E. Boser, I. M. Guyon, V. N. Vapnik, "A Training Algorithm for Optimal Margin Classifiers," in *Proc. of the Fifth Annual ACM Conference on Computational Learning Theory*, D. Haussler (eds.), ACM Press, Pittsburg, pp. 144-152, 1992.
- [9] L. Breiman, J. Olshen and C. Stone, "Classification and Regression Trees," Chapman and Hall, 1984.
- [10] P. Clarkson, and P. Moreno, "On The Use Of Support Vector Machines for Phonetic Classification," *Proc. of ICASSP 1999.*, pp. 585-588.
- [11] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and other kernel-based learning methods*, Cambridge University Press, 2000.
- [12] P. Comon, Independent component analysis, A new concept? *Signal Processing*, Vol. 36, pp. 287-314, 1994.
- [13] F. Cucker and S. Smale, "On the mathematical foundations of learning," *Bull. Am. Math. Soc.*, Vol. 39, pp. 1-49, 2002.
- [14] R. Duda, P. Hart and D. Stork, *Pattern Classification*, Wiley and Sons, New York, 2001.
- [15] FastICA Web Page, <http://www.cis.hut.fi/projects/ica/fastica/index.shtml>.
- [16] K. Fukunaga, *Statistical Pattern Recognition*, Academic Press, New York, 1989.
- [17] A. Ganapathiraju, J. Hamaker and J. Picone, "Support Vector Machines for Speech Recognition," *Proc. of ICSLP 2000*, Beijing, China, pp. 2931-2934.
- [18] M. G. Genton, "Classes of Kernels for Machine Learning: A Statistics Perspective," *J. Machine Learning Res.*, Vol. 2., pp. 299-312, 2001.
- [19] O. Ghitza, "Auditory Nerve Representation Criteria for Speech Analysis/Synthesis," *IEEE Trans. on ASSP*, Vol. 35., No. 6., pp. 736-740, 1987.
- [20] J. Glass, "A Probabilistic Framework for Segment-Based Speech Recognition," *Computer, Speech, and Language*, Vol. 17., pp. 137-152, 2003.
- [21] S. Greenberg, M. Slaney (eds.) "Computational Models of Auditory Function," IOS Press, 2001.
- [22] R. Haeb-Umbach and H. Ney, "Linear Discriminant Analysis for Improved Large Vocabulary Speech Recognition," *Proc. of ICASSP'92*, San Francisco, USA, March 1992, pp. 13-16.
- [23] H. Hermansky, "Modulation Spectrum in Speech Processing," In A. Prochazka et al. (eds.), *Signal Analysis and Prediction*, Birkhauser, 1998.
- [24] X. Huang, A. Acero and H.-W. Hon, *Spoken Language Processing*, Prentice Hall, 2001.
- [25] J.-N. Hwang, S.-R. Lay, M. Maechler, R. D. Martin and J. Schimert, Regression Modeling in Back-Propagation and Projection Pursuit Learning, *IEEE Trans. on Neural Networks*, Vol. 5., No. 3., pp. 342-353, 1994.
- [26] A. Hyvärinen, "A family of fixed-point algorithms for independent component analysis," *Proceedings of ICASSP*, Munich, Germany, 1997.
- [27] A. Hyvärinen, "New Approximations of Differential Entropy for Independent Component Analysis and Projection Pursuit," In *Advances in Neural Information Processing Systems*, Vol. 10, pp. 273-279, MIT Press, 1998.
- [28] A. Hyvärinen, "Fast and Robust Fixed-Point Algorithms for Independent Component Analysis," *IEEE Transactions on Neural Networks*, Vol. 10, No. 3, pp. 626-634, 1999.
- [29] I. J. Jolliffe, *Principal Component Analysis*, Springer-Verlag, New York, 1986.
- [30] B.-H. Juang and S. Katagiri, "Discriminative Learning for Minimum Error Classification". *IEEE Transactions on Signal Processing*, Vol. 40, No. 12, pp. 3043-3053., 1992.
- [31] S. Katagiri, C.-H. Lee and B.-H. Juang, "New Discriminative Training Algorithms Based on the Generalized Descent Method". *Proceedings of the 1991 IEEE Workshop on Neural Networks for Signal Processing*, pp. 299-308., 1991.
- [32] Kernel Machines Web site, <http://kernel-machines.org>.
- [33] A. Kocsor, L. Tóth, A. Kuba Jr., K. Kovács, M. Jelasity, T. Gyimóthy, J. Csirik, "A Comparative Study of Several Feature Transformation and Learning Methods for Phoneme Classification," *Int. Journal of Speech Technology*, Vol. 3., no. 3/4, pp. 263-276, 2000.
- [34] A. Kocsor, L. Tóth, D. Paczolay, "A Nonlinearized Discriminant Analysis and its Application to Speech Impediment Therapy," in V. Matousek et al. (eds.): *Proc. of Text, Speech and Dialogue 2001*, Springer Verlag LNAI Series, Vol. 2166, pp. 249-257, 2001.
- [35] A. Kocsor and J. Csirik, "Fast Independent Component Analysis in Kernel Feature Spaces," in: L. Pacholski and P. Ruzicka (eds.): *Proc. of SOFSEM 2001*, Springer Verlag LNCS Series, Vol. 2234, pp. 271-281, 2001.
- [36] A. Kocsor and K. Kovács, "Kernel Springy Discriminant Analysis and its Application to a Phonological Awareness Teaching System," in P. Sojka et al.(eds.): *Proc. of Text, Speech and Dialogue 2002*, Springer Verlag LNAI Series, Vol. 2448, pp. 325-328, 2002.
- [37] A. Kocsor and K. Kovács, unpublished result, 2003.
- [38] K.-F. Lee and H. Hon, Speaker-Independent Phone recognition Using Hidden Markov Models, *IEEE Trans. ASSP*, 37(11):1641-1648, November 1989.
- [39] A. Lima, H. Zen, Y. Nankaku, C. Miyajima, K. Tokuda and T. Kitamura, "On the Use of Kernel PCA for Feature Extraction in Speech Recognition," *Proc. of Eurospeech 2003*, Geneva, pp. 2625-2628.
- [40] S. Mika, G. Rätsch, J. Weston, B. Schölkopf, K.-R. Müller, "Fisher discriminant analysis with kernels," in: Y.-H. Hu et al. (eds.): *Neural Networks for Signal Processing IX*, pp. 41-48. *IEEE*, 1999.
- [41] B. C. J. Moore, *An Introduction to the Psychology of Hearing*, Academic Press, 1997.
- [42] Y. Normandin, "Hidden Markov Models, Maximum Mutual Information Estimation, and the Speech Recognition Problem", PhD Thesis, Department of Electrical Engineering, McGill University, Montreal., 1991.
- [43] P. Pudil, J. Novovicova, J. Kittler, "Floating search methods in feature selection," *Pattern Recognition Letters*, Vol. 15, pp. 1119-1125, 1994.
- [44] J. R. Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann, 1993.
- [45] R. Rosipal, L. J. Trejo, "Kernel Partial Least Squares Regression in Reproducing Kernel Hilbert Space," *J. Machine Learning Res.* Vol. 2, pp. 97-123, 2001.
- [46] J. Salomon, S. King and M. Osborne, "Framewise Phone Classification Using Support Vector Machines," *Proc. of ICSLP 2002*, pp. 2645-2648.
- [47] B. Schölkopf, A. J. Smola, K.-R. Müller, Kernel principal component analysis in *Advances in Kernel Methods - Support Vector Learning*, B. Schölkopf et al. (eds.), MIT Press, Cambridge, pp. 327-352, 1999.
- [48] B. Schölkopf, P. L. Bartlett, A. J. Smola, R. Williamson, Shrinking the tube: a new support vector regression algorithm in *Advances in Neural Information Processing Systems 11*, M. S. Kearns et al. (eds.), MIT Press, Cambridge, pp. 330-336, 1999.
- [49] B. Schölkopf and A. J. Smola, *Learning with Kernels*, MIT Press, Cambridge, 2002.
- [50] A. Shimodaira, K. Noma, M. Nakai, and S. Sagayama, "Support Vector Machine with Dynamic Time-Alignment Kernel for Speech Recognition," *Proc. of Eurospeech 2001*, Scandinavia, pp. 1841-1844.
- [51] O. Siohan, "On the Robustness of Linear Discriminant Analysis as a Preprocessing Step for Noisy Speech Recognition," *Proc. of ICASSP 1995*, Detroit, USA, May 1995, pp. 125-128.
- [52] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 1995.
- [53] V. N. Vapnik, *Statistical Learning Theory*, John Wiley & Sons Inc., 1998.
- [54] P.C. Woodland and D. Povey, "Large scale discriminative training of hidden Markov models for speech recognition," *Computer Speech and Language*, vol. 16, pp. 25-47, 2002.