

Second-Order Measures of Quality for Binary Classification: a Critical Overview and Their Use for Nonlinear Receiver Design

Fahed Abdallah

Cédric Richard

Régis Lengellé

Université de Technologie de Troyes (UTT), France

When deriving a detector, we are often led to consider design criteria such as second-order measures of quality. The aim of this paper is to provide a critical overview of these criteria. We first consider the case of deriving unconstrained detectors. We show that second-order criteria must satisfy a non-trivial condition to yield Bayes-optimal receivers, to be considered as relevant criteria for detector design. Next, we address the case where constraints are imposed on the detection structure, leading us to consider some set \mathcal{D} of admissible detectors. In these conditions we prove that even if there exists a monotonic function of the likelihood ratio in \mathcal{D} , obtaining this statistic via the optimization of a second-order criterion, relevant or not, is not guaranteed. Results are illustrated by simulation examples. Finally, in order to derive nonlinear discriminants via optimization of second-order criteria, we propose a method based on the kernel trick used in the implementation of the well-known support vector machine method. The new method is tested on a number of real data sets.

Keywords: Maximum likelihood, Detection, Distance measures, Signal-to-noise ratio, nonlinear discriminants.

INTRODUCTION

Let (\mathbf{X}, Y) be a pair of random variables taking their respective values from \mathbb{R}^d and $\{0, 1\}$, where \mathbf{X} is the observation and Y indicate either class \mathcal{C}_0 or \mathcal{C}_1 . The purpose of detection is to determine to which of two classes ($Y = 0$ or $Y = 1$) a given observation \mathbf{X} belongs. According to classical statistical detection theories, comparing any strictly monotonic function of the likelihood ratio $L(\mathbf{X})$ with a threshold value is the optimum test [13]. In practical applications, implementing such a test may be impossible because of incomplete specification of the conditional probability densities $p(\mathbf{X}|Y = 0)$ and $p(\mathbf{X}|Y = 1)$, denoted by the standard notations $p_0(\mathbf{X})$ and $p_1(\mathbf{X})$, respectively. Therefore we are often led to consider alternative design criteria such as second-order measures of quality. These criteria are easy to use since they only depend on first and second-order moments of the statistics S to be sought [6, 7]. A wide variety of second-order measures of performance have been proposed and several contributions have been presented to prove their efficiency, e.g., [8] and references therein.

In particular, some of these criteria guarantee the best solution in the Bayes sense since their optimization leads to a monotonic function of the likelihood ratio, as has been shown for well-known criteria such as Fisher criterion, mean-square error and signal-to-noise ratio.

The aim of this paper is to provide an overview of the strengths and shortcomings of second-order criteria. First, we consider the case of deriving unconstrained detectors. We show that second-order criteria must satisfy a non-trivial condition to provide Bayes-optimal receivers, to be considered as relevant second-order criteria for detector design. Next, we address the case where constraints are imposed on the structure of the detector, leading us to restrict our attention to some set \mathcal{D} of admissible detectors. In these conditions, we prove that even if a monotonic function of the likelihood ratio in \mathcal{D} exists, obtaining this statistic via the optimization of a relevant criterion is not guaranteed. Finally, an original method to derive nonlinear discriminants via optimization of second-order criteria is introduced. It uses the same kernel trick as the support vector machines (SVMs). It allows us to develop a nonlinear generalization of linear receivers obtained

via the optimization of second-order criteria. Simulation results are presented in order to compare the resulting nonlinear method with that of SVMs.

CHARACTERIZATION OF RELEVANT SECOND-ORDER CRITERIA

Background And Notations

Let $S(\mathbf{X}) : \mathbb{R}^d \rightarrow \mathbb{R}$ be an arbitrary measurable function and let $g : \mathbb{R}^d \rightarrow \{0, 1\}$ be the decision function based on $S(\mathbf{X})$:

$$g(\mathbf{X}) = \begin{cases} 1 & \text{if } S(\mathbf{X}) > b \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

which errs on \mathbf{X} if $g(\mathbf{X}) \neq Y$. Classical statistical theories such as Bayes, Neyman-Pearson and minimax lead to the fundamental result that the optimum test consists of comparing the likelihood ratio $L(\mathbf{X}) \triangleq p_1(\mathbf{X})/p_0(\mathbf{X})$ with a given threshold b in order to make a decision [13]. The decision rule can then be expressed as

$$g^*(\mathbf{X}) = \begin{cases} 1 & \text{if } L(\mathbf{X}) > b \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Note that g is equivalent¹ to the Bayes-optimal detector g^* , which means that their receiver operating characteristic (ROC) is the same, if $S(\mathbf{X}) = \phi\{L(\mathbf{X})\}$, where ϕ is any monotonic function. Since the implementation of (2) may be impossible in many practical applications, we are often led to consider simpler procedures for designing (1). In particular, one can use alternative design criteria such as second-order measures of performance. These criteria are defined in terms of first and second-order moments of the statistics $S(\mathbf{X})$, namely

$$m_i \triangleq E\{S | Y = i\}, \quad \sigma_i^2 \triangleq \text{Var}\{S | Y = i\}, \quad (3)$$

with $i \in \{0, 1\}$. There have been many contributions to justify the use of individual second-order criteria, including convergence with optimal detectors of classical detection theories (see, e.g., [8] and references therein). Examples of these criteria are Fisher, mean-square error and signal-to-noise ratio. In [7, pp. 141-3], the objective of the author is to unify these results stating that the use of any function $\Psi(m_0, m_1, \sigma_0^2, \sigma_1^2)$ as a criterion for general non-linear detector design leads to a Bayes-optimum detector. In fact, we show in the next subsection that $\Psi(m_0, m_1, \sigma_0^2, \sigma_1^2)$ must satisfy a non-trivial condition to guarantee Bayes-optimal detectors for general nonlinear detector design, and

¹Throughout this paper, two detectors are said to be equivalent if their receiver operating characteristic are the same.

thus to be considered as a *relevant second-order criterion* for detector design.

Relevant Second-Order Criteria

Let Ψ be any function of $m_i \triangleq \int S(\mathbf{X}) p_i(\mathbf{X}) d\mathbf{X}$ and $\sigma_i^2 \triangleq \int (S(\mathbf{X}) - m_i)^2 p_i(\mathbf{X}) d\mathbf{X}$, where $S(\mathbf{X})$ denotes any decision statistic. We first have to characterize statistics $S(\mathbf{X})$ which optimize Ψ . Operating on Ψ with a variational calculus, we obtain

$$\delta\Psi = \frac{\partial\Psi}{\partial m_0} \delta m_0 + \frac{\partial\Psi}{\partial m_1} \delta m_1 + \frac{\partial\Psi}{\partial \sigma_0^2} \delta \sigma_0^2 + \frac{\partial\Psi}{\partial \sigma_1^2} \delta \sigma_1^2. \quad (4)$$

Since $\delta m_i = \int \delta S(\mathbf{X}) p_i(\mathbf{X}) d\mathbf{X}$ and $\delta \sigma_i^2 = 2 \int (S(\mathbf{X}) - m_i) \delta S(\mathbf{X}) p_i(\mathbf{X}) d\mathbf{X}$ with $i \in \{0, 1\}$, we obtain

$$\begin{aligned} \delta\Psi = \int \left[\frac{\partial\Psi}{\partial m_0} p_0(\mathbf{X}) + \frac{\partial\Psi}{\partial m_1} p_1(\mathbf{X}) \right. \\ \left. + 2(S(\mathbf{X}) - m_0) \frac{\partial\Psi}{\partial \sigma_0^2} p_0(\mathbf{X}) \right. \\ \left. + 2(S(\mathbf{X}) - m_1) \frac{\partial\Psi}{\partial \sigma_1^2} p_1(\mathbf{X}) \right] \delta S(\mathbf{X}) d\mathbf{X}. \end{aligned} \quad (5)$$

In order to make $\delta\Psi = 0$ regardless of $\delta S(\mathbf{X})$, the $[\cdot]$ term in the integrand must be equal to 0. Using $L(\mathbf{X}) = \frac{p_1(\mathbf{X})}{p_0(\mathbf{X})}$, we finally get the expression of the statistic $S(\mathbf{X})$ optimizing Ψ as a function of the likelihood ratio

$$S(\mathbf{X}) = -\frac{1}{2} \frac{\frac{\partial\Psi}{\partial m_0} + \frac{\partial\Psi}{\partial m_1} L(\mathbf{X})}{\frac{\partial\Psi}{\partial \sigma_0^2} + \frac{\partial\Psi}{\partial \sigma_1^2} L(\mathbf{X})} + \frac{m_0 \frac{\partial\Psi}{\partial \sigma_0^2} + m_1 \frac{\partial\Psi}{\partial \sigma_1^2} L(\mathbf{X})}{\frac{\partial\Psi}{\partial \sigma_0^2} + \frac{\partial\Psi}{\partial \sigma_1^2} L(\mathbf{X})} \quad (6)$$

The above statistic $S(\mathbf{X})$ leads to a Bayes-optimal detector if, and only if, it is a strictly monotonic function of $L(\mathbf{X})$. Evaluating the first order derivative of $S(\mathbf{X})$ with respect to $L(\mathbf{X})$, we obtain

$$\frac{dS}{dL}(\mathbf{X}) = \frac{(m_1 - m_0) \frac{\partial\Psi}{\partial \sigma_0^2} \frac{\partial\Psi}{\partial \sigma_1^2} + \frac{1}{2} \left(\frac{\partial\Psi}{\partial \sigma_1^2} \frac{\partial\Psi}{\partial m_0} - \frac{\partial\Psi}{\partial \sigma_0^2} \frac{\partial\Psi}{\partial m_1} \right)}{\left(\frac{\partial\Psi}{\partial \sigma_0^2} + \frac{\partial\Psi}{\partial \sigma_1^2} L(\mathbf{X}) \right)^2} \quad (7)$$

We thus note that $S(\mathbf{X})$ defined by (6) is a strictly monotonic function of $L(\mathbf{X})$ if, and only if, the numerator of (7) is not equal to 0. This result leads directly to the following proposition [14].

Proposition 1. $\Psi(m_0, m_1, \sigma_0^2, \sigma_1^2)$ is a relevant second-order criterion Ψ_R , i.e., it guarantees the best solution in the Bayes sense, if and only if

$$(m_1 - m_0) \frac{\partial\Psi}{\partial \sigma_0^2} \frac{\partial\Psi}{\partial \sigma_1^2} + \frac{1}{2} \left(\frac{\partial\Psi}{\partial \sigma_1^2} \frac{\partial\Psi}{\partial m_0} - \frac{\partial\Psi}{\partial \sigma_0^2} \frac{\partial\Psi}{\partial m_1} \right) \neq 0. \quad (8)$$

Since it is very difficult, if not impossible, to find the solutions of (8), the above property can only be used to test the relevance of any criteria $\Psi(m_0, m_1, \sigma_0^2, \sigma_1^2)$. However, note that (8) is a non-restrictive condition, i.e., there exists a broad class of second-order criteria that lead to Bayes-optimal detectors [1].

CONSTRAINED DETECTOR DESIGN USING SECOND-ORDER CRITERIA

As shown in the previous section, there exists a broad class of second-order criteria that lead to a detection statistic $S(\mathbf{X})$ equivalent to the likelihood ratio $L(\mathbf{X})$. However, implementing $S(\mathbf{X})$ remains an unsolved problem since it depends on the probability densities $p_0(\mathbf{X})$ and $p_1(\mathbf{X})$ via $L(\mathbf{X})$, which are unknown. Therefore, we have to use the following strategy for deriving receivers [5]:

1. selecting a class \mathcal{D} of detection statistics;
2. choosing the statistic of \mathcal{D} that optimizes a given measure of performance, e.g., a relevant second-order criterion.

Unfortunately, this approach does not necessarily provide a Bayes-optimal detector since it generally requires the optimum statistic (6) to be a member of \mathcal{D} . We shall now discuss this drawback in the case where \mathcal{D} denotes the class of linear statistics.

Linear Detectors Design

Linear detectors are the simplest to use as far as their implementation is concerned, and are directly related to many known techniques such as correlations and Euclidean distances [7]. We shall now show how second-order criteria can be used for designing linear detectors, which are defined as follows:

$$g(\mathbf{X}) = \begin{cases} 1 & \text{if } S(\mathbf{X}) = \mathbf{W}^T \mathbf{X} - b > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Here \mathbf{W} denotes the direction onto which any n -dimensional observation \mathbf{X} is projected, and b is the detector threshold. The conditional expected values and variances of $S(\mathbf{X})$ are given by

$$m_i = E\{S | Y = i\} = \mathbf{W}^T \mathbf{M}_i - b \quad (10)$$

$$\sigma_i^2 = \text{Var}\{S | Y = i\} = \mathbf{W}^T \boldsymbol{\Sigma}_i \mathbf{W}, \quad (11)$$

where \mathbf{M}_i and $\boldsymbol{\Sigma}_i$ are the conditional expected vectors and covariance matrices of \mathbf{X} . Let Ψ be any second-order criterion. The optimal statistic $S(\mathbf{X})$ is given by equating to zero the partial derivatives of Ψ with respect to \mathbf{W} and b . As shown in [7, pp. 133-4] and [12], solving this linear system leads directly to the following proposition.

Proposition 2. Let $S(\mathbf{X}) \triangleq \mathbf{W}^T \mathbf{X} - b$ be any linear decision statistic. The optimum projection vector \mathbf{W} under which the maximum value of any second-order criteria Ψ is reached satisfies

$$\mathbf{W}_\rho = [\rho \boldsymbol{\Sigma}_0 + (1 - \rho) \boldsymbol{\Sigma}_1]^{-1} [\mathbf{M}_1 - \mathbf{M}_0], \quad (12)$$

where \mathbf{M}_i and $\boldsymbol{\Sigma}_i$ are the conditional expected vectors and covariance matrices of \mathbf{X} . The parameter ρ depends on the criterion Ψ as follows:

$$\rho = \frac{\frac{\partial \Psi}{\partial \sigma_0^2}}{\frac{\partial \Psi}{\partial \sigma_0^2} + \frac{\partial \Psi}{\partial \sigma_1^2}}. \quad (13)$$

The optimum projection direction \mathbf{W}_ρ depends on Ψ through a single parameter $\rho \in] -\infty, +\infty[$. In this case, the latter can be chosen to optimize the performance of the detector. Note that $\rho \in [0, 1]$ if, and only if, $\partial \Psi / \partial \sigma_0^2$ and $\partial \Psi / \partial \sigma_1^2$ are of the same sign (**Property 1**). This condition means that Ψ varies in the same way with σ_0^2 and σ_1^2 , which is a desirable but non-mandatory requirement for design criteria. Let us now concentrate on $\mathbf{W}_{-\infty}$ and $\mathbf{W}_{+\infty}$. They are both proportional to $[\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1]^{-1} [\mathbf{M}_1 - \mathbf{M}_0]$ since we have

$$\mathbf{W}_{\pm\infty} \propto \lim_{\rho \rightarrow \pm\infty} \frac{1}{\rho} [\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1]^{-1} [\mathbf{M}_1 - \mathbf{M}_0]. \quad (14)$$

The projection directions $\mathbf{W}_{-\infty}$ and $\mathbf{W}_{+\infty}$ then lead to equivalent detection structures (**Property 2**) since the ROC depends only on the direction of \mathbf{W} .

In the following, Proposition 2 and the above properties are illustrated through some classic detection problems. As mentioned at the very beginning of this section, we also show that (relevant) second-order criteria does not guarantee an optimal detector in the Bayes sense if we restrict the solution space to a specific class \mathcal{D} of detectors. Here this drawback is illustrated through the following situations:

Scenario 1: the optimization of any second-order criteria in \mathcal{D} leads to a Bayes-optimal detector,

Scenario 2: certain second-order criteria exist which provide Bayes-optimal receivers in \mathcal{D} ,

Scenario 3: there exists a Bayes-optimal detector in \mathcal{D} but it cannot be reached by optimizing any second-order criteria,

where \mathcal{D} denotes the class of linear detectors (9).

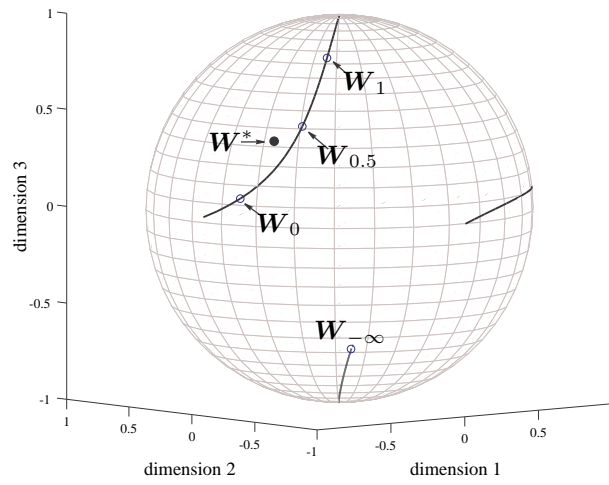


Fig. 1. Evolution of the projection direction \mathbf{W}_ρ on the unit sphere as a function of $\rho \in]-\infty, +\infty[$, in the case of three-dimensional exponential distributions ($\lambda_{01} = 5$, $\lambda_{11} = 2$, $\lambda_{02} = 3$, $\lambda_{12} = 2$, $\lambda_{03} = 2$, $\lambda_{13} = 3$). The projection direction associated with the Bayes-optimal detector is referred to as \mathbf{W}^* in this figure.

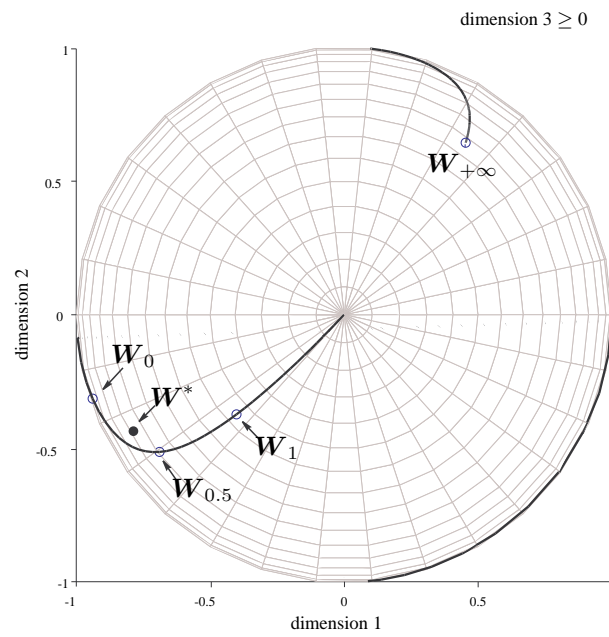


Fig. 2. Same as figure 1. Overhead view of the unit sphere, i.e., dimension 3 ≥ 0 .

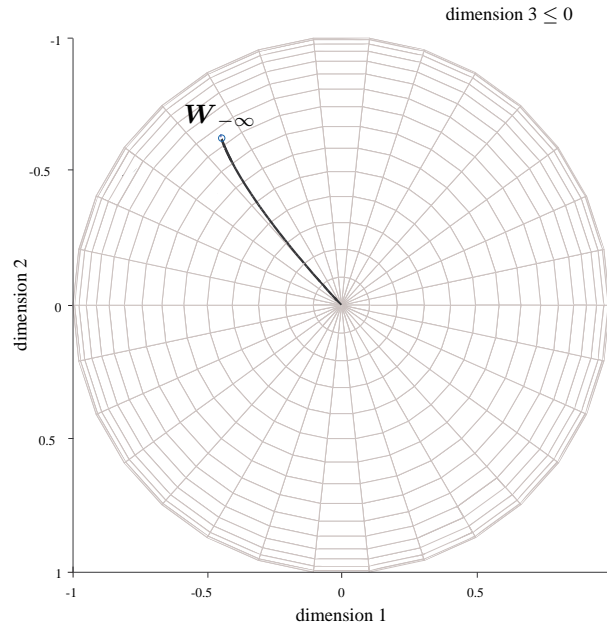


Fig. 3. Same as figure 1. View from below of the unit sphere, i.e., dimension 3 \leq 0.

Data set	KFD	SVM	KSOD
Banana	10.60	10.43 (132)	10.59
Thyroid	0.39	0.33 (156)	0.25
B. cancer	8.14	7.10 (364)	6.70
Diabetes	17.79	17.68 (308)	17.39
German	21.36	21.06 (400)	20.96
Heart	4.44	4.52 (388)	4.41
Solar	32.42	32.73 (364)	31.61
Waveform	11.14	11.07 (400)	11.14
Ringnorm	1.53	1.52 (396)	1.53
Titanic	28.88	28.88 (400)	28.55

Table 1. Comparison of the mean error rates obtained with KFD, SVM and KSOD. This table also gives the number of support vectors used by the SVM method.

Scenario 1: Case Of Normal Distributions With Equal Covariances

When $p_0(\mathbf{X})$ and $p_1(\mathbf{X})$ are normal with expected vectors \mathbf{M}_0 and \mathbf{M}_1 and covariance matrices Σ_0 and Σ_1 , it is well-known that the Bayes-optimal statistic is given by:

$$S(\mathbf{X}) = \frac{1}{2}(\mathbf{X} - \mathbf{M}_0)^T \Sigma_0^{-1}(\mathbf{X} - \mathbf{M}_0) - \frac{1}{2}(\mathbf{X} - \mathbf{M}_1)^T \Sigma_1^{-1}(\mathbf{X} - \mathbf{M}_1). \quad (15)$$

This equation shows that the decision boundary is a quadratic form in \mathbf{X} . When $\Sigma_0 = \Sigma_1 = \Sigma$, the boundary becomes a linear function of \mathbf{X} as

$$S(\mathbf{X}) = (\mathbf{M}_1 - \mathbf{M}_0)^T \Sigma^{-1} \mathbf{X}. \quad (16)$$

Eq. (16) indicates that the direction onto which any \mathbf{X} is projected is given by $\mathbf{W}^* = \Sigma^{-1}(\mathbf{M}_1 - \mathbf{M}_0)$. Let us now determine the projection direction under which the optimal value of any second-order criterion is reached. Eq. (12) gives:

$$\mathbf{W}_\rho = \Sigma^{-1}(\mathbf{M}_1 - \mathbf{M}_0). \quad (17)$$

Comparing the projection directions \mathbf{W}^* and \mathbf{W}_ρ , we immediately conclude that they both correspond to equivalent detection structures. This very simple example shows that any second-order criterion, relevant or not, can sometimes lead to a Bayes-optimal receiver. However, such a success is not always guaranteed as illustrated in the next subsection, even if it exists a Bayes-equivalent receiver in \mathcal{D} .

Scenarios 2 And 3: Case Of Exponential Distributions

Let us consider that the components X_j of \mathbf{X} are exponentially distributed and mutually independent. Then we have:

$$p_i(\mathbf{X}) = \prod_{j=1}^d \frac{1}{\lambda_{ij}} \exp\left(-\frac{1}{\lambda_{ij}} X_j\right) u(X_j), \quad i \in \{0, 1\}, \quad (18)$$

with λ_{ij} the parameter of the exponential distribution of the random variable X_j , and $u(\cdot)$ the step function. It can easily be shown that the linear function given below is the Bayes-optimal detection statistic:

$$S(\mathbf{X}) = \sum_{j=1}^d \left(\frac{1}{\lambda_{0j}} - \frac{1}{\lambda_{1j}} \right) X_j. \quad (19)$$

Then it is associated with the following Bayes-optimal projection direction:

$$\mathbf{W}^* = \left(\frac{\lambda_{11} - \lambda_{01}}{\lambda_{11}\lambda_{01}}, \dots, \frac{\lambda_{1j} - \lambda_{0j}}{\lambda_{1j}\lambda_{0j}}, \dots, \frac{\lambda_{1d} - \lambda_{0d}}{\lambda_{1d}\lambda_{0d}} \right). \quad (20)$$

The expected vector \mathbf{M}_i and the covariance matrix Σ_i of \mathbf{X} , which is exponentially distributed according to (18), are given by $\mathbf{M}_i = (\lambda_{i1}, \dots, \lambda_{ij}, \dots, \lambda_{id})^T$ and $\Sigma_i = \text{diag}(\lambda_{i1}^2, \dots, \lambda_{ij}^2, \dots, \lambda_{id}^2)$. Applying (12) to determine the projection direction under which the optimal value of any second-order criterion is reached, we obtain:

$$\mathbf{W}_\rho = \left(\frac{\lambda_{11} - \lambda_{01}}{\lambda_{11}^2 - \rho(\lambda_{11}^2 - \lambda_{01}^2)}, \dots, \frac{\lambda_{1j} - \lambda_{0j}}{\lambda_{1j}^2 - \rho(\lambda_{1j}^2 - \lambda_{0j}^2)}, \dots, \frac{\lambda_{1d} - \lambda_{0d}}{\lambda_{1d}^2 - \rho(\lambda_{1d}^2 - \lambda_{0d}^2)} \right). \quad (21)$$

Comparing (20) and (21) shows that the collinearity of \mathbf{W}^* and \mathbf{W}_ρ depends on ρ . We shall now illustrate Scenarios 2 and 3.

Consider the case of two-dimensional observations \mathbf{X} with $\lambda_{11} \neq \lambda_{01}$ and $\lambda_{12} \neq \lambda_{02}$. Vectors \mathbf{W}^* and \mathbf{W}_ρ are collinear if, and only if,

$$\rho = \frac{\lambda_{02}\lambda_{12}\lambda_{11}^2 - \lambda_{01}\lambda_{11}\lambda_{12}^2}{\lambda_{02}\lambda_{12}(\lambda_{11}^2 - \lambda_{01}^2) - \lambda_{01}\lambda_{11}(\lambda_{12}^2 - \lambda_{02}^2)}. \quad (22)$$

This means that any second-order criterion Ψ guarantees the best solution in the Bayes sense if, and only if, it satisfies:

$$\frac{\frac{\partial \Psi}{\partial \sigma_0^2}}{\frac{\partial \Psi}{\partial \sigma_0^2} + \frac{\partial \Psi}{\partial \sigma_1^2}} = \frac{\lambda_{02}\lambda_{12}\lambda_{11}^2 - \lambda_{01}\lambda_{11}\lambda_{12}^2}{\lambda_{02}\lambda_{12}(\lambda_{11}^2 - \lambda_{01}^2) - \lambda_{01}\lambda_{11}(\lambda_{12}^2 - \lambda_{02}^2)}, \quad (23)$$

an example of which is the generalized signal-to-noise ratio Ψ_α :

$$\Psi_\alpha(S) = \frac{(m_1 - m_0)^2}{(1 - \alpha)\sigma_1^2 + \alpha\sigma_0^2}, \quad (24)$$

the parameter α must then be equal to the second term in (22). Here m_i and σ_i^2 denote the conditional expected values and variances of S . This illustrates Scenario 2.

Consider now that \mathbf{X} is a d -dimensional observation, with ($d > 2$). Except for very particular cases, we notice that there is no ρ parameter which ensures the collinearity of \mathbf{W}^* and \mathbf{W}_ρ . This means that optimization of second-order criteria, relevant or not, does not necessarily lead to a Bayes-equivalent detector, even if there is one that is a member of \mathcal{D} . Figures 1, 2 and 3 illustrate this situation, called Scenario 3, for three-dimensional observations \mathbf{X} . The projection direction \mathbf{W}_ρ is represented on the unit sphere as a function of $\rho \in]-\infty, +\infty[$. One can observe the collinearity of projection directions $\mathbf{W}_{-\infty}$ and $\mathbf{W}_{+\infty}$ (see Property 2). Finally, one can notice that there is no value ρ_0

such that \mathbf{W}_{ρ_0} and \mathbf{W}^* are collinear, which is one of the weaknesses of second-order criteria.

NONLINEAR SECOND-ORDER DISCRIMINANT

In recent years a great interest has been shown in kernel-based algorithms for developing a nonlinear generalization of linear receivers, see [11] and references therein. Kernel-based classification algorithms were primarily used in Support Vector Machines (SVMs) [3, 4]. By mapping the samples $(\mathbf{X}_i)_{i=1,\dots,n}$ into a high dimensional feature space and reformulating the problem into dot product form in order to use Mercer kernels, an effective solution for nonlinear discriminant analysis has been obtained [17, chapter 5]. This exploits the notion that applying a nonlinear data transformation to some high-dimensional feature space increases the probability of having linearly separable classes in the transformed space. In [10], a nonlinear classification technique based on Fisher discriminants has been proposed. It also uses the Mercer kernel trick, and allows the efficient computation of linear Fisher discriminants in feature space. Very promising results had been reported using this approach, called the kernel Fisher discriminant method (KFD), when compared with other state of the art classification techniques. In this paper, we present an extension of the KFD method that also deals with nonlinear discriminant analysis using kernel functions and second-order measures of performance. This method is based on extending the expression (12) to the nonlinear case by mapping the data by a nonlinear transformation.

Formulation Of The Nonlinear Discriminant

Applying a nonlinear data transformation to some specific high dimensional feature spaces increases the probability of having linearly separable classes within the transformed space [17, 11]. Linear discriminants in the feature space are then equivalent to nonlinear discriminants in the original space. More formally, by mapping the samples $\{\mathbf{X}_i\}_{i=1,\dots,n}$ using a nonlinear function

$$\begin{aligned} \Phi: \mathbb{R}^d &\longrightarrow \mathcal{F} \\ \mathbf{X} &\longmapsto \Phi(\mathbf{X}), \end{aligned}$$

one can perform a linear discriminant analysis in \mathcal{F} with the set $\{(\Phi(\mathbf{X}_i), Y_i)\}_{i=1,\dots,n}$ in order to obtain a nonlinear discriminants in the original space.

Clearly, if \mathcal{F} is a very high, or even infinitely, dimensional space, deriving $S(\mathbf{X}) = (\mathbf{W})^T \Phi(\mathbf{X}) - b$ may be a computationally intractable problem. However, by using the theory of reproducing kernels [3], such a problem can be solved without explicitly mapping the data to the feature space \mathcal{F} .

Nonlinear kernel second-order discriminant (KSOD) can be obtained by using (12) in the feature space \mathcal{F} . According to Proposition 1, the function $S(\mathbf{X})$ operating in \mathcal{F} is optimum in the sense of any given second-order criterion Ψ if it satisfies

$$[\rho \Sigma_0^\Phi + (1 - \rho) \Sigma_1^\Phi] \mathbf{W} = [\mathbf{M}_1^\Phi - \mathbf{M}_0^\Phi], \quad (25)$$

where \mathbf{M}_i^Φ and Σ_i^Φ denote the conditional expected vectors and covariance matrices of $\Phi(\mathbf{X})$, respectively. These moments can be estimated as follows:

$$\mathbf{M}_i^\Phi = \frac{1}{n_i} \sum_{\mathbf{X} \in \mathcal{C}_i} \Phi(\mathbf{X}) \quad (26)$$

$$\Sigma_i^\Phi = \frac{1}{n_i} \sum_{\mathbf{X} \in \mathcal{C}_i} \Phi(\mathbf{X}) \Phi^T(\mathbf{X}) - (\mathbf{M}_i^\Phi) (\mathbf{M}_i^\Phi)^T, \quad (27)$$

where n_i is the number of samples from class \mathcal{C}_i in the training set. When \mathcal{F} is a very high-dimensional space, (25) may be difficult to solve except if the Mercer trick is used. From the theory of reproducing kernels [15], we know that any solution $\mathbf{W} \in \mathcal{F}$ must lie in the span of all training samples in \mathcal{F} . Therefore \mathbf{W} can be written as follows:

$$\mathbf{W} = \sum_{i=1}^n \alpha(i) \Phi(\mathbf{X}_i) = \mathbf{Q} \boldsymbol{\alpha} \quad (28)$$

where \mathbf{Q} denotes the matrix $[\Phi(\mathbf{X}_1) \cdots \Phi(\mathbf{X}_n)]$, and the $\alpha(i)$'s are the dual parameters. Multiplying (25) by \mathbf{Q}^T and using (28) yields

$$[\rho \mathbf{Q}^T \Sigma_0^\Phi \mathbf{Q} + (1 - \rho) \mathbf{Q}^T \Sigma_1^\Phi \mathbf{Q}] \boldsymbol{\alpha} = \mathbf{Q}^T [\mathbf{M}_1^\Phi - \mathbf{M}_0^\Phi]. \quad (29)$$

Let k be any kernel that satisfies the Mercer condition [17, 11]. Then we have

$$k(\mathbf{X}_i, \mathbf{X}_j) = \Phi(\mathbf{X}_i)^T \Phi(\mathbf{X}_j), \quad (30)$$

which means that $k(\mathbf{X}_i, \mathbf{X}_j)$ corresponds to the inner product of \mathbf{X}_i and \mathbf{X}_j in \mathcal{F} .

By using a kernel which verifies (30), the expression (29) can be reformulated as

$$\mathbf{P}_\rho \boldsymbol{\alpha} = \mathbf{M}, \quad (31)$$

where $\boldsymbol{\alpha}$ has to be determined and \mathbf{P}_ρ is a n by n matrix which is given by

$$\mathbf{P}_\rho = \left[\frac{\rho}{n_0} \mathbf{K}_0 (\mathbf{I} - \mathbf{1}_{n_0}) \mathbf{K}_0^T + \frac{1 - \rho}{n_1} \mathbf{K}_1 (\mathbf{I} - \mathbf{1}_{n_1}) \mathbf{K}_1^T \right]. \quad (32)$$

In the above expression, \mathbf{K}_i is a n by n_i matrix with elements

$$\mathbf{K}_i(p, q) = k(\mathbf{X}_p, \mathbf{X}_q), \quad (33)$$

for all $\mathbf{X}_p \in (\mathcal{C}_0 \cup \mathcal{C}_1)$ and $\mathbf{X}_q \in \mathcal{C}_i$. \mathbf{I} is the identity matrix and $\mathbf{1}_{n_i}$ is the matrix with all elements set to $\frac{1}{n_i}$. The components of M in (31) are defined as

$$M(j) = \frac{1}{n_1} \sum_{\mathbf{X} \in \mathcal{C}_1} k(\mathbf{X}, \mathbf{X}_j) - \frac{1}{n_0} \sum_{\mathbf{X} \in \mathcal{C}_0} k(\mathbf{X}, \mathbf{X}_j). \quad (34)$$

To determine the projection of any new sample \mathbf{X} onto \mathbf{W} , we have to calculate the n -dimensional vector α from (31). Eq. (28) yields:

$$\Phi(\mathbf{X})^T \mathbf{W} = \sum_{i=1}^n \alpha(i) k(\mathbf{X}_i, \mathbf{X}). \quad (35)$$

Finally one can use different strategies to determine the bias b , as shown in [10]. We can summarize our work in the algorithm described below to determine a minimum error rate detector from a class of kernel-based decision structures that are optimum in the sense of second-order criteria.

1. Given any Mercer kernel k , compute the kernel matrices \mathbf{K}_0 and \mathbf{K}_1 from (33).
2. Compute M from (34).
3. Set ρ to zero.
4. While ($\rho \leq 1$) repeat
 - Compute \mathbf{P}_ρ from (32).
 - Solve (31) to get the vector α , and retain the result as α_ρ .
 - Find the threshold b_ρ which minimizes, e.g., an estimate of the generalization error.²
 - Update ρ : $\rho \leftarrow \rho + \Delta\rho$, where $\Delta\rho$ is a selected step.
5. Select the best detector characterized by (α_ρ, b_ρ) .

From equation (13), one can see that the parameter ρ varies from 0 to 1 if the derivatives of the corresponding second-order criteria Ψ with respect to σ_0^2 and σ_1^2 have the same sign. This desirable but non-mandatory requirement for design criteria makes the process much simpler than adjusting ρ in \mathbb{R} .

The method presented above is called *kernel second-order discriminant* since it leads to the optimum nonlinear receiver in the sense of the best second-order criterion without setting it up. Obviously, classifiers obtained with KSOD perform better than or

equal to those resulting from the KFD method developed by Mika *et al.* in [10].

Experimentations

To compare our method to KFD and SVM, 10 experiments were conducted on artificial and real world data downloaded from <http://www.first.gmd.de/~raetsch>. For each of the 10 problems, Table 1 shows the average test error over 40 runs on 400 training samples and 8000 test samples chosen arbitrarily from a mixture of the available data sets. The kernel function was selected as the RBF having a width equal to 1 [17]. Note that the numerical problems created by inverting the ill-conditioned matrix \mathbf{P} in (31) can be avoided by adding a regularization term. It can be chosen as a multiple of the identity matrix [10, 11, 17], i.e., replace \mathbf{P}_ρ by the matrix $\mathbf{P}_\rho + \eta\mathbf{I}$ with $\eta > 0$. The results presented in Table 1 clearly show that the KSOD method can often perform favourably compared with the other state of the art detection techniques. However, note that for the KSOD method, all the training data should be used to test a new sample. This is not the case for the SVM method which uses only the training samples called support vectors [17]. Hence, the testing time for the KSOD method is in general higher than that of the SVMs. Table 1 gives the number of support vectors used by the SVM method.

CONCLUSION

The theoretical results reported in this paper are concerned with the virtues and vices of second-order criteria used for detector design. Firstly, we have given a necessary and sufficient condition for these measures of performance to guarantee the best solution in the Bayes sense when deriving unconstrained detectors. Secondly we have considered the case where constraints are imposed on the structure of detectors, leading us to restrict our attention to a class \mathcal{D} of admissible detectors. We have shown that any second-order criterion, relevant or not, can sometimes lead to a Bayes-optimal receiver. However, such a success is far from assured in the majority of cases, even if a Bayes-equivalent detector in \mathcal{D} exists.

Finally, we have proposed a nonlinear classification technique based on second-order criteria. It uses the Mercer kernel trick, and allows for the efficient computation of linear second-order discriminants in feature space. The results obtained suggest that the method presented in this paper often performs better than SVM and KFD. However, one should note that the testing time for this method is in general higher than that of the SVMs since it uses all the training samples in the test phase. Future works will be dedicated to control the complexity of the resulting discriminant

²See, e.g., [10] for other criteria and strategies.

in order to increment the generalization performances and to reduce the testing time [11, 16, 9].

REFERENCES

- [1] F. Abdallah, C. Richard and R. Lengell. On equivalence between detectors obtained from second-order measures of performance. *XI European Signal Processing Conference*, September 3-6, 2002, Toulouse, France.
- [2] S. A. Billings and K. L. Lee. "Nonlinear Fisher discriminant analysis using a minimum squared error cost function and the orthogonal least squares algorithm," *Neural Networks*, vol. 15, pp. 263-270, 2002.
- [3] B. Boser, I. Guyon and V. Vapnik. "A training algorithm for optimal margin classifiers," *Proc. Fifth Annual Workshop on Computational Learning Theory*, Pittsburgh, pp. 144-152, 1992.
- [4] C. Cortes and V. Vapnik. "Support vector networks," *Machine learning*, vol. 20, pp. 273-297, 1995.
- [5] L. Devroye, L. Györfi and G. Lugosi. *A probabilistic theory of pattern recognition*, Springer-Verlag, 1996.
- [6] R. O. Duda, P. E. Hart and D. C. Stork. *Pattern classification*. New York: John Wiley and Sons, 2001.
- [7] K. Fukunaga. *Statistical Pattern Recognition*. Academic Press, 1990.
- [8] W. A. Gardner. A unifying view of second-order measures of quality for signal classification. *IEEE Transactions on Communications*, vol. 28, no. 6, pp. 807-816, 1980.
- [9] Y. LeCun, J. S. Denker and S. A. Solla. "Optimal brain damage," in *Proc. Advances in Neural Information Processing Systems*, vol. 2, pp. 598-605, 1990.
- [10] S. Mika, G. Rätsch, J. Weston, B. Schölkopf and K. R. Müller. Fisher discriminant analysis with kernels, in *Advances in Neural Networks for Signal Processing*, Y. H. Hu, J. Larsen, E. Wilson, S. Douglas, editors, pp. 41-48, 1999.
- [11] K. R. Müller, S. Mika, G. Rätsch, K. Tsuda and B. Schölkopf. "An introduction to kernel-based learning algorithms," *IEEE Neural Networks*, vol. 12, no. 2, pp. 181-201, May 2001.
- [12] D. W. Patterson and R. L. Mattson. A method of finding linear discriminant functions for a class of performance criteria. *IEEE Transactions on Information Theory*, vol. 12, no. 3, pp. 380-387, 1966.
- [13] H. V. Poor. *An introduction to signal detection and estimation*. Springer-Verlag, 1994.
- [14] C. Richard, R. Lengell and F. Abdallah. Bayes-optimal detectors design using relevant second-order criteria. *IEEE Signal Processing Letters*, vol. 9, no. 1, 2002.
- [15] S. Saitoh. *Theory of Reproducing Kernels and its Applications*. Longman Scientific & Technical, 1988.
- [16] J. Rissanen. "Stochastic complexity in statistical inquiry," *Annals of Statistics*, vol. 14, no. 3, pp. 1080-1100, 1986.
- [17] V. Vapnik. *The nature of statistical learning theory*. New York: Springer Verlag, 1995.

NOTATIONS

All vectors and matrices are in bold

\mathcal{D}	class of functions or detectors
C_i	class i
\mathbf{X}	input sample
$X(i)$	the component number i of \mathbf{X}
Y_i	target value for class i
d	dimensionality
n	number of training samples
\mathcal{F}	high dimensional feature space
$p_i(\mathbf{X})$	conditional probability density of class i
$S(\cdot)$	measurable function
$g(\cdot)$	decision function
$L(\cdot)$	likelihood ratio
Ψ	second order criterion
m_i	expected value with respect to class i
σ_i	variance with respect to class i
\mathbf{W}	d -dimensional vector
α	n dimensional dual vector
$\alpha(i)$	the component number i of α
$k(\cdot, \cdot)$	Mercer kernel
Φ	nonlinear function
M_i	conditional expected vector of class i
M_i^Φ	conditional expected vector of class i in feature space
Σ_i	conditional covariance matrix of class i
Σ_i^Φ	conditional covariance matrix of class i in feature space
ROC	Receiver Operating Characteristic
KFD	Kernel Fisher Detectors
SVM	Support Vector Machine
KSOD	Kernel Second Order Detectors