



# A novel formulation of orthogonal polynomial kernel functions for SVM classifiers: The Gegenbauer family

Luis Carlos Padierna<sup>a</sup>, Martín Carpio<sup>a</sup>, Alfonso Rojas-Domínguez<sup>a,1,\*</sup>, Héctor Puga<sup>a</sup>, Héctor Fraire<sup>b</sup>

<sup>a</sup> Tecnológico Nacional de México, Instituto Tecnológico de León, León 37290, México

<sup>b</sup> Tecnológico Nacional de México, Instituto Tecnológico de Cd. Madero, Cd. Madero 89460, México

## ARTICLE INFO

### Article history:

Available online 12 July 2018

### Keywords:

SVM classifier  
Orthogonal polynomials  
Gegenbauer kernel  
Binary classification

## ABSTRACT

Orthogonal polynomial kernels have been recently introduced to enhance support vector machine classifiers by reducing their number of support vectors. Previous works have studied these kernels as isolated cases and discussed only particular aspects. In this paper, a novel formulation of orthogonal polynomial kernels that includes and improves previous proposals (Legendre, Chebyshev and Hermite) is presented. Two undesired effects that must be avoided in order to use orthogonal polynomial kernels are identified and resolved: the Annihilation and the Explosion effects. The proposed formulation is studied by means of introducing a new family of orthogonal polynomial kernels based on Gegenbauer polynomials and comparing it against other kernels. Experimental results reveal that the Gegenbauer family competes with the RBF kernel in accuracy while requiring fewer support vectors and overcomes other classical and orthogonal kernels.

© 2018 Elsevier Ltd. All rights reserved.

## 1. Introduction

Support Vector Machines (SVMs) represent a machine learning model used to solve pattern recognition tasks such as classification, regression and outlier detection [1]. SVMs belong to a family of models for pattern analysis known as kernel-based learning methods that also include support vector clustering, kernel principal components analysis, kernel perceptron, etc. [2]; these methods are motivated by rigorous theoretical analysis and characterized by computational efficiency. Kernel methods enable practitioners to analyze nonlinear relations in data via the use of the so-called kernel trick. The first and arguably the most popular kernel method to date is the SVM. In the context of pattern classification, SVMs were originally developed to handle two-class problems [3,4]. Their solution involves a quadratic programming problem, with dual form [5]:

$$\begin{aligned} \max L(\lambda) &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \lambda_i \lambda_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t. } C \geq \lambda_i \geq 0 \quad \forall i = 1, \dots, N, \text{ and } \sum_{i=1}^N \lambda_i y_i &= 0 \end{aligned} \quad (1)$$

where  $C$  is a scalar parameter called the penalty factor,  $\lambda$  are non-negative Lagrange multipliers [4], and  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  represent a set of training data;  $\mathbf{x}_i \in \mathbb{R}^d$  is the  $i$ -th input vector and  $y_i \in \{+1, -1\}$  its corresponding class label. Ultimately, those data points  $\mathbf{x}_i$  for which the corresponding multipliers  $\lambda_i \neq 0$  are the support vectors that define the decision boundary (a hyperplane in the feature space  $\mathcal{H}$ , which in general is a nonlinear boundary in the original space).

The function  $K(\mathbf{x}, \mathbf{z})$  in (1), defined on  $\mathbb{R}^d \times \mathbb{R}^d$  is called a kernel if there exists a mapping to the feature (Hilbert) space  $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$  such that  $K(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$  [4]. The performance of SVMs is highly dependent on the kernel function, which can take different forms [6]. A variety of kernel functions have been developed for different problem domains, e.g. Linear [7]; String [8]; non-parametric [9]; etc. The difference between these kernels is in their own ability to represent the similarities between the data points in the feature space. Kernel optimization, kernel selection, and kernel design, have thus become important tasks in machine learning. In this paper, we address the problem of kernel design through the study of orthogonal polynomial kernels, and introduce a new family of kernels based on Gegenbauer polynomials.

In the last two decades, orthogonal polynomials (OPs) have attracted interest as valid kernels for SVMs [6,10–14]. The most relevant studies are discussed in Section 2. The main idea behind building orthogonal polynomial kernels relies on the fact that the mapping  $\phi(x)$  for  $x \in \mathbb{R}$  can be defined

\* Corresponding author.

E-mail address: [alfonso.rojas@gmail.com](mailto:alfonso.rojas@gmail.com) (A. Rojas-Domínguez).

<sup>1</sup> CONACYT research fellow.

in terms of  $n$ -th degree polynomials  $\{p_n(x)\}_{n=0}^{\infty}$  as  $\phi(x) := (p_0(x), p_1(x), p_2(x), \dots, p_n(x))$ . Previous works on the generation of orthogonal polynomial kernels have shown that the number of support vectors in the SVM tends to be smaller when OPs are used as kernels than when classical kernel functions are used [13–15]. According to the theory developed by Vapnik in 1998 [6], the expectation of the probability of error for the optimal hyperplane constructed on the basis of training samples of size  $N$  (the expectation taken over both training and test data) is upper-bounded according to:

$$\mathbb{E}R(\lambda^{(N)}) \leq \frac{\mathbb{E} \min \left\{ \mathcal{K}^{(N+1)}, \left( \frac{\mathcal{D}^{(N+1)}}{\rho^{(N+1)}} \right)^2 \right\}}{N+1} \quad (2)$$

where  $\mathcal{K}$  is the number of (essential) support vectors,  $\mathcal{D}$  is the maximum norm among support vectors,  $\rho$  is the maximal possible margin between the separating hyperplane and its closest vector and the superscript notation indicates the size of the training dataset. Thus, it follows that selecting the solution with the smallest number of support vectors (and/or increasing the amount of training data) can improve the generalization of an SVM (cf. Chapter 10 of [6]).

The main hypothesis that has been put forward to explain the relatively fewer support vectors required when using OPs as kernels is that these help to decrease data redundancy in the feature space. Our own conjecture is that orthogonal polynomial kernels induce kernel matrices with fewer significant eigenvalues and thus, less support vectors. Recently, generalization measures have been proposed based on spectral analysis of the kernel matrix [16] that may support these hypotheses. The goal in designing the optimal kernel function is therefore to achieve high classification performance with the smallest possible number of support vectors [17]. Optimality shall thus be measured under two criteria: the classification accuracy and the Proportion of Support Vectors (PSV) used to define the decision boundary:  $\text{PSV} = 100 \times \mathcal{K}/N$ .

The main contributions of this paper are: (i) Previous proposals of orthogonal polynomial kernels are analyzed and positioned within a novel formulation that satisfies the Mercer conditions [18] while correcting some problems encountered within previous formulations. (ii) The hypothesis that better orthogonal polynomial kernels can be constructed with the so far unexplored members of the ultraspherical-polynomials is put forward and, to test this hypothesis, a new Gegenbauer kernel family is introduced. Only one extra control parameter is required to select a specific member of this family. (iii) A comprehensive experimental evaluation of Gegenbauer against classical and orthogonal polynomial kernels is carried out through a robust methodology that emphasizes the optimization of hyper-parameters, so that kernels are compared on the fairest ground possible. It is shown that the proposed Gegenbauer kernel overcomes its predecessors in terms of the considered indexes: Classification Accuracy and PSV.

The remaining sections of this paper are structured as follows: Section 2 provides a brief theoretical background of orthogonal polynomials and a review of previous orthogonal polynomial kernels. In Section 3, our novel formulation to produce orthogonal polynomial kernels for classification is described and the Gegenbauer kernel functions constructed under this formulation are presented. The experimental methodology is described in Section 4, and experimental results are reported and discussed in Section 5. Finally, conclusions to this work and further directions to extend this research are offered in Section 6.

## 2. Related work

### 2.1. Orthogonal polynomials

Orthogonal polynomials (OPs) can be regarded as universal function-approximators used in a variety of fields; their main quality is to reduce the redundancy in data, and they have succeeded in applications such as interpolation, data compression, quadrature approximation [19] and recently as pattern recognition techniques for feature extraction [20] and regression [21]. A complete survey on the theory supporting OPs can be consulted in [22]. However, the following definitions summarize the theory needed to explain the construction of orthogonal polynomial kernels.

**Orthogonal Polynomial (OP).** A polynomial  $p(x)$  is said to be an orthogonal polynomial with respect to an inner product  $\langle \cdot, \cdot \rangle$  if, given the space of  $n$ -th degree polynomials of one variable,  $\Pi_n$ , the inner product [23]:

$$\langle p(x), q(x) \rangle = 0 \quad \forall p(x) \in \Pi_n, \quad \text{degree}(q(x)) \neq \text{degree}(p(x)) \quad (3)$$

**Orthogonal Polynomial Sequence (OPS).** If the inner product  $\langle p_m(x), p_n(x) \rangle$  is defined via a non-negative weight function  $w(x)$  integrable on the open interval  $(a, b)$ , the sequence of  $n$ -th degree polynomials of one variable  $\{p_n(x)\}_{n=0}^{\infty}$  is called an orthogonal polynomial sequence with respect to  $w(x)$  on  $(a, b)$  [24]; and the condition of orthogonality is given by:

$$\langle p_m(x), p_n(x) \rangle = \int_a^b p_m(x) p_n(x) w(x) dx = 0, \quad \forall m \neq n \quad (4)$$

If instead of the continuous interval,  $(a, b)$ ,  $w(x)$  is defined on a set of isolated points  $V$  in  $\mathbb{R}$ ; then  $\{p_n(x)\}_{n=0}^{\infty}$  is an OPS  $\{p_n: 0 \leq n \leq |V|\}$ , that satisfies [25]:

$$\langle p_m(x), p_n(x) \rangle = \sum_{x \in V} p_m(x) p_n(x) w(x) = 0, \quad \forall m \neq n \quad (5)$$

The condition of orthogonality entails the property that any inner product between a pair of orthogonal polynomials of different degrees becomes zero. The condition is satisfied by polynomials from every known family of orthogonal polynomials, such as: Legendre, Chebyshev, Gegenbauer (ultraspherical), Laguerre, Hermite and Jacobi (characteristics of all of these polynomials can be consulted in the Askey scheme, cf. chapter 18 of [26]).

The property just described can be directly employed in the formulation of kernels based on orthogonal polynomials by including only products between polynomials of the same degree in the kernel design; these kernels are called Orthogonal Polynomial Kernels. Previous proposals for construction of this type of kernels are reviewed in the next subsection.

### 2.2. Previous orthogonal polynomial kernels

In 1998, Vapnik [6] described a number of kernels, including kernels that generate expansion of polynomials; in particular, he described a kernel based on Hermite polynomials  $H(x)$  as follows:

$$K_H(x, z) = \sum_{i=0}^{\infty} q^i H_i(x) H_i(z) \quad (6)$$

where  $0 \leq q \leq 1$  is a convergence factor. Vapnik [6] also demonstrated that, for multidimensional basis functions that are tensor products of the coordinatewise basis, a  $d$ -dimensional kernel is obtained by the product of  $d$  kernels of the form in (6):

$$K(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d K_j(x_j, z_j) \quad (7)$$

In 2006, Ye et al. reported some of the first practical experiments with orthogonal polynomial kernels, and proposed a Chebyshev kernel in terms of univariate polynomials as [10]:

$$K_{Cheby}(x, z) = \frac{\sum_{i=0}^n T_i(x)T_i(z)}{\sqrt{1-xz}} \quad (8)$$

where  $(\sqrt{1-xz})^{-1}$  corresponds to the weight function  $w(x, z)$  and  $T_i(x)$  is the  $i$ -th degree Chebyshev type-I polynomial. In order to compute the kernel function from input vectors  $\mathbf{x}$  and  $\mathbf{z}$  of dimension  $d$ , the following equation is used [13]:

$$K_{Cheby}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \frac{\sum_{i=0}^n T_i(x_j)T_i(z_j)}{\sqrt{1-x_jz_j}} \quad (9)$$

Similarly, other orthogonal polynomial kernels can be produced by substituting each  $T_i(x_j)$  by the desired polynomials and an appropriate weight function. For instance, the Legendre kernel can be formulated using the Legendre polynomials  $P_i(x_j)$  with  $w(x_j, z_j) = 1$  as [11]:

$$K_{Legen}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n P_i(x_j)P_i(z_j) \quad (10)$$

In 2008, Ozer and Chen [12] extended the work of Ye et al. by expressing the kernel function in vector form originally, instead of doing it for each scalar, and defined the Generalized Chebyshev Kernel (GCK) and the Modified Chebyshev Kernel (MCK):

$$K_{GCK}(\mathbf{x}, \mathbf{z}) = \frac{\sum_{i=0}^n T_i(\mathbf{x})T_i^T(\mathbf{z})}{\sqrt{d - \langle \mathbf{x}, \mathbf{z} \rangle}} \quad (11)$$

$$K_{MCK}(\mathbf{x}, \mathbf{z}) = \frac{\sum_{i=0}^n T_i(\mathbf{x})T_i^T(\mathbf{z})}{\exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)} \quad (12)$$

where  $T_i(\mathbf{x})$  is a row vector. In 2011, Ozer et al. experimented with GCK and MCK on few instances of pattern classification problems [13]. A more comprehensive set of experiments on pattern classification using SVMs with RBF kernel and GCK (among other classifiers) were carried out recently by Sun et al. [15]. In spite of the fact that the PSV was not used to measure the performance of the SVMs, the classification accuracies reported on twenty one test datasets provide an important reference for our experimental comparison.

In 2016, Moghaddam and Hamidzadeh showed that the use of modified Hermite ( $He(x)$ ) polynomials in the formulation of a kernel function does indeed lead to a decrease in the PSV [14]. They defined the Hermite-based kernel as:

$$K_{Herm}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n He_i(x_j)He_i(z_j) \quad (13)$$

Moghaddam and Hamidzadeh evaluated their proposed kernel on fifteen real-world datasets and conclude that it also improves classification accuracy and training time [14].

The orthogonal polynomial kernels discussed in this work are listed in Table 1. The *Annihilation* and *Explosion* effects referred to in this table are described in later sections. A more detailed analysis of each individual kernel can be found in the corresponding references provided. In this work, our aim is to perform a robust and uniform evaluation of these kernels (with the exception of the one proposed by Vapnik for which there are no experimental results reported in the literature) and compare them against our own proposal.

Our proposal follows a construction method similar to the formulation first described by Vapnik and used in Eqs. (7), (9), (10) and (13). The difference between previous proposals and our formulation lies in the fact that we include a weight function of

two variables and a scaling function designed to eliminate undesired effects that may hinder the operation of the orthogonal polynomial kernels. This is described in detail in the next section, where new kernels based on the Gegenbauer family of polynomials are developed following our novel formulation.

### 3. A novel formulation of orthogonal polynomial kernels

This section describes our proposal of a new formulation of orthogonal polynomial SVM kernels. In this formulation, an SVM kernel is constructed as the tensor product of the inner product of weighted and scaled univariate polynomials, i.e. kernel functions of scalar inputs which are then aggregated (under Mercer closures) to produce a vector-form kernel. Denoting the weight function by  $w(x_j, z_j)$  and the scaling function by  $u(p_i)$ , the general form of an orthogonal polynomial SVM kernel under our formulation is thus given by:

$$K(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle = \prod_{j=1}^d \sum_{i=0}^n p_i(x_j)p_i(z_j)w(x_j, z_j)u(p_i)^2 \quad (14)$$

The weight function is independent of the degree of the polynomials; this function is not unique, but it is selected from known functions to satisfy the orthogonality condition for each particular family of polynomials. The scaling function,  $u(p_i) \equiv \beta |p_i(\cdot)|_{\max}^{-1} \in \mathbb{R}^+$ , is the reciprocal of the maximum absolute value of an  $i$ -th degree polynomial within its operational range, and  $\beta$  is a convenient positive scalar factor (see Section 3.1.2). Eq. (14) is general in the sense that it allows the creation of new kernels and it encompasses most previous proposals.

It has been observed [10,13,21] that a kernel combination of the form  $\prod_{j=1}^d K(x_j, z_j)$  will yield a very small value if even a single one  $K(x_j, z_j)$  is close to zero. Hereinafter this undesired effect is referred to as the *Annihilation Effect*. In trying to overcome this problem, some researchers in previous works [10,13,21] have empirically modified the weight function in the denominator of each  $K(x_j, z_j)$  by adding a small constant  $\varepsilon \ll 1$ . For instance, the denominator  $\sqrt{1-x_jz_j}$  in (9) becomes  $\sqrt{1-x_jz_j+\varepsilon}$  in experiments [27]. We propose a different way to deal with the annihilation effect. Recall that for every family of orthogonal polynomials,  $p_0(x) = 1$ ; thus, it immediately follows that the premise  $K(x_j, z_j) \rightarrow 0$  can only occur if  $\sum_{i=1}^n p_i(x) \rightarrow -1$ . The proposed Eq. (14) includes a scaling function  $u(p_i)$  that modifies the series  $\sum_{i=1}^n p_i(x)$  such that each  $|p_i(x)| \leq 1$ . Details will be provided in Section 3.1.2.

A second undesired effect occurs whenever  $K(x_j, z_j) \gg 1$  and/or the dimensionality  $d$  of the problem is relatively big. In this case the total output of  $|\prod_{j=1}^d K(x_j, z_j)| \rightarrow \infty$ . This effect, deemed the *Explosion Effect*, will lead to numerical problems in the implementation of the kernels. The proposed scaling function in (14) is designed to avoid this problem. Details and proofs are contained in the sections below.

In the next section a new family of orthogonal polynomials is presented. The use of this family in the construction of orthogonal polynomial kernels poses a number of advantages, such as: increased capability to deal with high-dimensional problems, analytical tractability, and inclusion of other kernels.

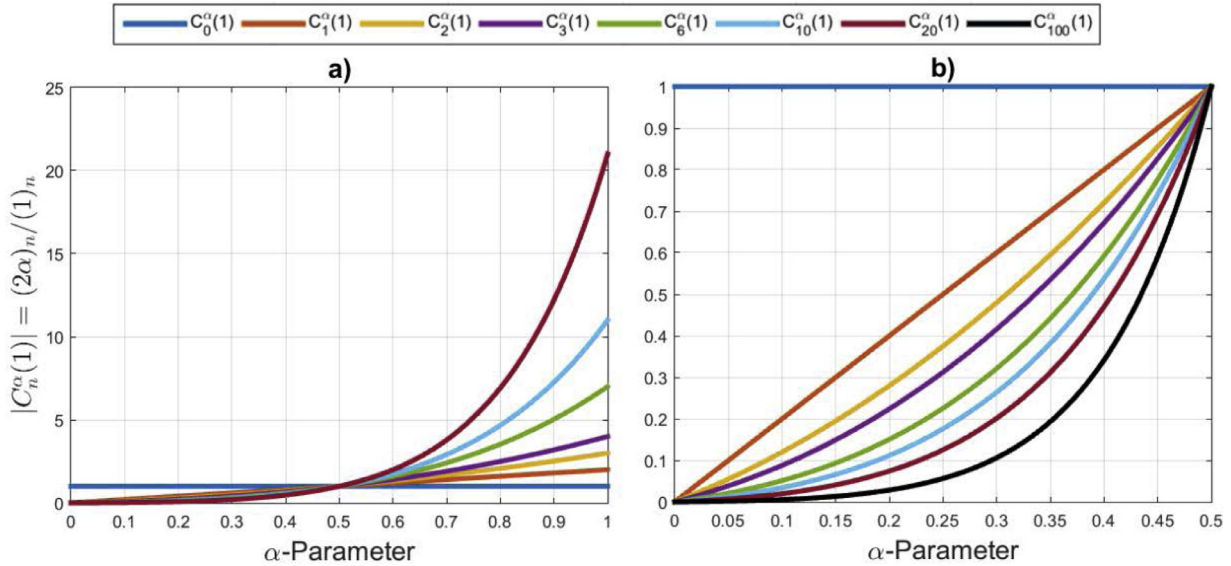
#### 3.1. The Gegenbauer kernel functions

A new family of kernels that follows our proposed formulation established in Eq. (14) is presented. This family is based on Gegenbauer polynomials, which includes both Chebyshev and Legendre kernels. The following discussion describes the family of polynomials as well as its corresponding weight and scaling functions.

**Table 1**

List of the orthogonal kernels discussed in this work.

Kernel (main ref.)	Short Label	Observations
Hermite [6]	$K_H$	First OP kernel described. Involves an infinite series with a convergence factor. No experimental results have been found on this kernel.
Chebyshev [10]	$K_{Cheby}$	Reports the first experimental results with OP kernels for SVM classifiers. Suffers from the <i>Annihilation</i> and <i>Explosion</i> effects.
Generalized Chebyshev [13]	$K_{GCK}$	First OP kernel formulated directly in vector form. Improves $K_{Cheby}$ and defines one way to avoid the <i>Annihilation</i> effect.
Modified Chebyshev [13]	$K_{MCK}$	Kernel formulated in vector form (including an exponential function). Has shown better performance on non-linear problems than $K_{GCK}$ .
Legendre [11]	$K_{Legen}$	Its advantage is that it does not require a weight function. It is free from the <i>Explosion</i> effect, but suffers from the <i>Annihilation</i> effect.
Hermite [14]	$K_{Herm}$	Prone to the <i>Annihilation</i> and <i>Explosion</i> effects. Hermite polynomials are orthogonal in $(-\infty, \infty)$ , requiring manipulation of the input data. No weight function was originally defined for this kernel.
This Work, Eq. (20)	$K_{Gegen}$	Constructed according to Eq. (14). Avoids the <i>Annihilation</i> and <i>Explosion</i> effects. Includes and improves $K_{Cheby}$ and $K_{Legen}$ as special cases.
This Work, Eq. (30)	$K_{S-Herm}$	Constructed according to Eq. (14). Avoids the <i>Annihilation</i> and <i>Explosion</i> effects. Corrects problems of the $K_{Herm}$ kernel.

**Fig. 1.** Maximum values obtained by the scaling function for the Gegenbauer polynomials. a) Case of  $\alpha \in (0, 1]$ ; b) case of  $\alpha \in (0, 0.5]$  where the amplitude grows with  $\alpha$  but decreases with the degree.

### 3.1.1. Gegenbauer polynomials

The Gegenbauer polynomials of degree  $n$  and parameter  $\alpha$ , denoted as  $C_n^\alpha(x)$ , are solutions of the Gegenbauer differential equation [28]:  $(1 - x^2)y'' - (2\alpha + 1)xy' + n(n + 2\alpha)y = 0$  and produced by means of the three-term recurrence equation, Eq. (15) with  $C_0^\alpha(x) = 1$ ,  $C_1^\alpha(x) = \alpha 2x$ . An infinite number of ultraspherical polynomials can be derived from the Gegenbauer family, one for each possible value of  $\alpha > -0.5$ . These polynomials along with an appropriate weight function  $w_\alpha(x)$  satisfy (4) for  $-1 \leq x \leq 1$ ; i.e., the generated polynomials are orthogonal in this range of  $x$  with respect to  $w_\alpha(x)$  [29]. By substituting  $\alpha = 0.5$ , and  $\alpha = 1$  in Eq. (15), the Legendre and Chebyshev type-II polynomials are obtained, respectively. Any other value of  $\alpha > -0.5$  can be used, with the exception of  $\alpha = 0$  (Chebyshev type-I), in whose case Eq. (16) should be used. Fig. 1 shows the maximum amplitudes of the Gegenbauer polynomials as function of  $\alpha$ .

$$(n + 1)C_{n+1}^\alpha(x) = 2(n + \alpha)x C_n^\alpha(x) - (n + 2\alpha - 1)C_{n-1}^\alpha(x) \quad (15)$$

$$C_{n+1}^{\alpha=0}(x) = 2xC_n^{\alpha=0}(x) - C_{n-1}^{\alpha=0}(x); \text{ with } C_0^{\alpha=0}(x) = 1, C_1^{\alpha=0}(x) = x \quad (16)$$

### 3.1.2. The weight and scaling functions

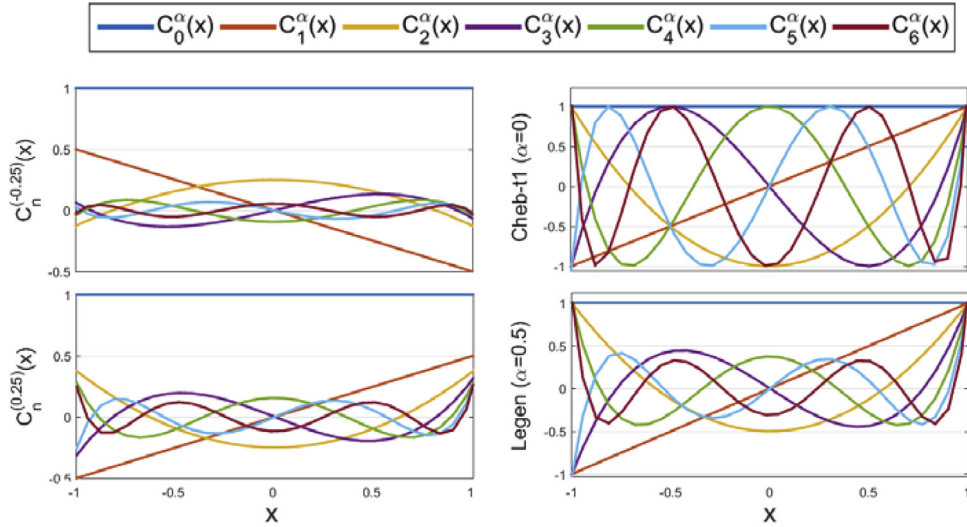
The Gegenbauer polynomials can be classified into two groups according to the parameter  $\alpha$ : the first group for  $-0.5 < \alpha \leq 0.5$  (including two special cases: Chebyshev type-I for  $\alpha = 0$  and Legendre for  $\alpha = 0.5$ ), and the second group with  $\alpha > 0.5$ . Regarding the first group, observe in Fig. 2 that their amplitude  $|C_n^\alpha(1)| \leq 1$ , such that a scaling function is not required for polynomials in this group. Similarly, the weight function for this group equals one.

In the second group, the growth of the polynomials amplitude to values  $\gg 1$  is evident (Fig. 3a illustrates two cases), which can lead to the Explosion Effect. A mechanism to avoid this undesired effect is given in the form of the weight and scaling functions. The weight function is fixed for each family of orthogonal polynomials; in the case of Gegenbauer this is [23,29]:

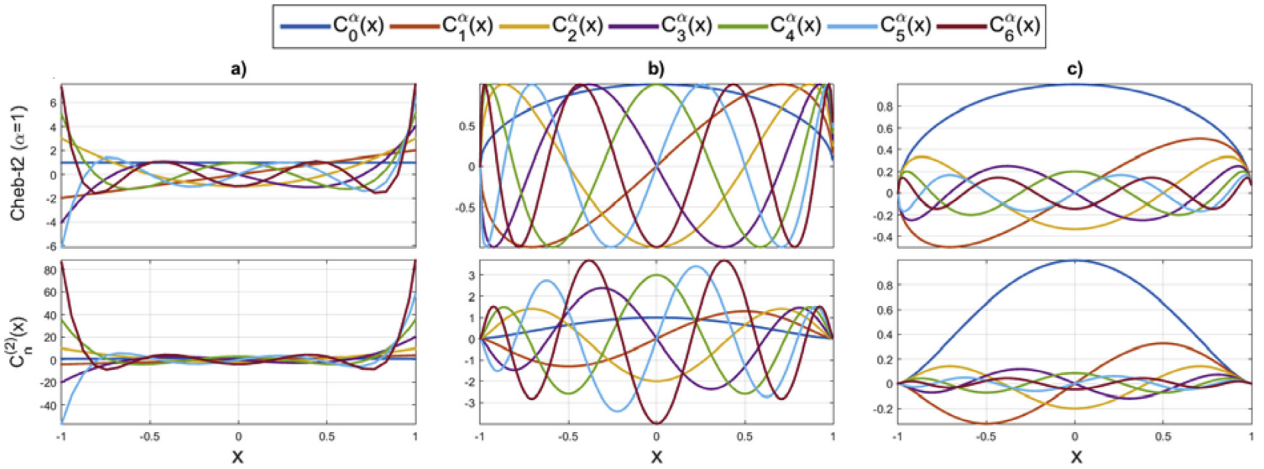
$$w_\alpha(x) = (1 - x^2)^{\alpha-1/2} \quad (17)$$

This function was used to produce the polynomials in Fig. 3b. Still, for  $\alpha > 1$  it can be observed that the magnitude of the polynomials  $\gg 1$ . Thus, to completely eliminate the Explosion and Annihilation effects, the Gegenbauer polynomials in the second group need to not only be weighted but also scaled (Fig. 3c). According to our formulation, the scaling function is  $u(p_i) = \beta |p_i(\cdot)|_{\max}^{-1} \in \mathbb{R}^+$ . The Gegenbauer polynomials possess the property that their maximum amplitudes are always reached for  $x = \pm 1$ . Consequently, the





**Fig. 2.** Examples of Gegenbauer polynomials in the first group. These polynomials do not require a scaling function since their amplitude  $|C_n^\alpha(1)| \leq 1$ .



**Fig. 3.** a) Original, b) weighted, c) weighted-and-scaled Gegenbauer polynomials in the second group. These polynomials require weighting and scaling to eliminate the Annihilation and Explosion effects.

scaling function can be computed by means of the Pochhammer operator  $(a \in \mathbb{R})_{n \in \mathbb{Z}} = a(a+1)(a+2) \cdots (a+n-1)$ , given  $\alpha > 0.5$  and  $i > 0$  as:  $|C_i^\alpha(1)| = (2\alpha)_i / (i!)_i$ . Setting  $\beta = 1/\sqrt{n+1}$ , the scaling function becomes:

$$u(p_i) = u(C_i^\alpha) = \left( \sqrt{n+1} |C_i^\alpha(1)| \right)^{-1} \quad (18)$$

### 3.1.3. Weight function of two variables

The analysis presented above is valid for univariate polynomials; however, although the shape of the weight function is preserved and independent of input data, when using Gegenbauer polynomials as kernels, two variables must be considered. Following the recent theory developed by Dunkl and Xu [23], the bivariate weight function is defined by taking the product of the corresponding univariate functions:

$$w_\alpha(x, z) = ((1-x^2)(1-z^2))^{\alpha-1/2} + \varepsilon \quad (19)$$

For the first group of Gegenbauer polynomials  $(-0.5 < \alpha \leq 0.5)$ , the weight function is approximately 1 for every point except nearby the extremes, as shown in Fig. 4a. Further, the amplitude of these polynomials is already small enough not to require the scaling by a weight function. Consequently, for this group of polynomials a convenient approximation is to use  $w_\alpha(x, z) = 1$  which

substantially simplifies the computation of the corresponding kernels.

For the second group of polynomials  $(\alpha > 0.5)$ , the weight function becomes zero at values on the border (Fig. 4b), which causes problems, particularly if a dataset contains categorical features (because many instances could be mapped to these extreme values causing data loss). In this case, adding a small offset (e.g.  $\varepsilon = 0.1$ ) solves the problem.

Up to this point, we have described all the elements required to formulate the Gegenbauer kernel functions. Now we can define the Gegenbauer kernel as:

$$K_{\text{Geg}}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n C_i^\alpha(x_j) C_i^\alpha(z_j) w_\alpha(x_j, z_j) u(C_i^\alpha)^2 \quad (20)$$

### 3.2. Theoretical properties

A valid kernel function must satisfy the necessary and sufficient conditions established in the Mercer's theorem [18]. Briefly stated, this theorem affirms that, given a finite input space  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ , the series  $\sum_{i=1}^\infty v_i \theta_i(\mathbf{x}) \theta_i(\mathbf{z})$ , in terms of eigenfunctions  $\theta_i \in L_2(\mathbf{X})$  and positive associated eigenvalues  $v_i$ , converges absolutely and uniformly to  $K(\mathbf{x}, \mathbf{z})$  when the latter is symmetric

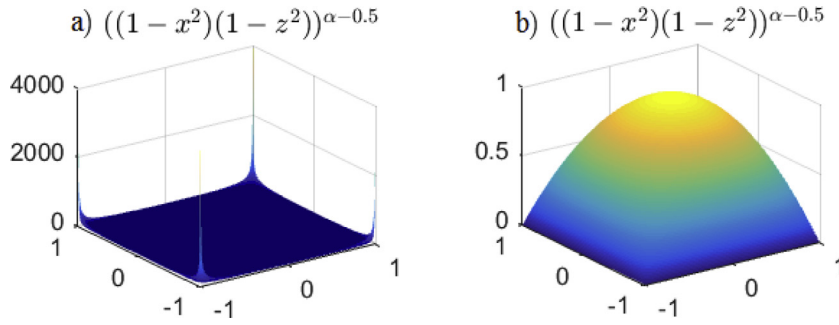


Fig. 4. Weight functions for the a) first group and b) second group of Gegenbauer kernels.

and positive semidefinite. Combinations of these kernels such as: summation or multiplication between two kernels; multiplication of a kernel with a positive real number; or composition of two kernels, produce other valid kernels because these combinations are defined as closures. For a more formal statement of this theorem and proofs of the closures, cf. Section 3.3 in [30]. Below, three important theoretical properties of the Gegenbauer kernel are described.

### 3.2.1. The Gegenbauer kernel is a valid Mercer kernel

A valid Mercer kernel is symmetric and positive semidefinite (non-negative definite). The proof that (20) is symmetric is trivial. We shall focus on demonstrating that it is positive semidefinite. According to the Mercer's theorem [18]: given a symmetric and continuous function  $K(x, z)$  defined in the closed square  $a \leq x \leq b$ ,  $a \leq z \leq b$ ; and any function  $g(\cdot)$  that ranges the class of all functions which are continuous in the closed interval  $[a, b]$ , the sufficient condition for  $K(x, z)$  to be positive is:

$$\int_a^b \int_a^b K(x, z) g(x) g(z) dx dz \geq 0 \quad (21)$$

**Theorem 1.** The function given in Eq. (20) is a positive Mercer kernel.

**Proof.** first notice that  $K_{Geg}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d K_{Geg}(x_j, z_j)$ , so that the Gegenbauer kernel defined for scalar data  $x$  and  $z$  is expressed as:

$$K_{Geg}(x, z) = \sum_{i=0}^n C_i^\alpha(x) C_i^\alpha(z) w_\alpha(x, z) u(C_i^\alpha)^2 \quad (22)$$

Substitution of (22) into (21) leads to (the integral limits are ignored for compactness):

$$\begin{aligned} & \iint K_{Geg}(x, z) g(x) g(z) dx dz \\ &= \iint \sum_{i=0}^n C_i^\alpha(x) C_i^\alpha(z) w_\alpha(x, z) u(C_i^\alpha)^2 g(x) g(z) dx dz \end{aligned}$$

substituting (19), which is separable:

$$\begin{aligned} &= \iint \sum_{i=0}^n C_i^\alpha(x) C_i^\alpha(z) \left[ (1-x^2)^{\alpha-1/2} (1-z^2)^{\alpha-1/2} + \varepsilon \right] \\ &\quad \times u(C_i^\alpha)^2 g(x) g(z) dx dz \end{aligned}$$

by definition,  $u(C_i^\alpha)$  is always positive and independent of the data, hence:

$$\begin{aligned} &= \sum_{i=0}^n u(C_i^\alpha)^2 \iint C_i^\alpha(x) C_i^\alpha(z) (1-x^2)^{\alpha-1/2} (1-z^2)^{\alpha-1/2} g(x) g(z) dx dz \\ &\quad + \sum_{i=0}^n \varepsilon u(C_i^\alpha)^2 \iint C_i^\alpha(x) C_i^\alpha(z) g(x) g(z) dx dz \end{aligned}$$

$$\begin{aligned} &= \sum_{i=0}^n u(C_i^\alpha)^2 \int C_i^\alpha(x) (1-x^2)^{\alpha-1/2} g(x) dx \int C_i^\alpha(z) (1-z^2)^{\alpha-1/2} g(z) dz \\ &\quad + \sum_{i=0}^n \varepsilon u(C_i^\alpha)^2 \int C_i^\alpha(x) g(x) dx \int C_i^\alpha(z) g(z) dz \\ &= \sum_{i=0}^n u(C_i^\alpha)^2 \left( \int C_i^\alpha(x) (1-x^2)^{\alpha-1/2} g(x) dx \right)^2 \\ &\quad + \sum_{i=0}^n \varepsilon u(C_i^\alpha)^2 \left( \int C_i^\alpha(x) g(x) dx \right)^2 \geq 0 \quad (23) \end{aligned}$$

Thus,  $K_{Geg}(x, z)$  is positive semidefinite, and it has been proved that it is a valid Mercer kernel.

Secondly, as the product of two kernels is also a kernel [1],  $K_{Geg}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d K_{Geg}(x_j, z_j)$  for vector data  $\mathbf{x}$  and  $\mathbf{z}$  is also a valid Mercer kernel under this closure. ■

### 3.2.2. The Gegenbauer kernel prevents the adverse annihilation effect

The annihilation effect occurs when the summation in  $K(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n p_i(x_j) p_i(z_j)$  produces 0 for at least one of the  $d$  dimensions. In the case of the Gegenbauer kernel this effect can only occur if  $\sum_{i=0}^n C_i^\alpha(x_j) C_i^\alpha(z_j) w_\alpha(x_j, z_j) u(C_i^\alpha)^2 = -1$ , since  $C_0^\alpha(x) C_0^\alpha(z) w_\alpha(x, z) |C_0^\alpha(1)|^{-2} = 1$ . Therefore, to prove that this effect is avoided, it must be proved that the summation for degrees  $i = 1 : n$  cannot be equal to  $-1$ .

**Proof.** for any  $-1 < x < 1$ ,  $-1 < z < 1$ ,  $n \geq 1$ , and  $\alpha > -0.5$ :

$$\left| \sum_{i=1}^n C_i^\alpha(x) C_i^\alpha(z) w_\alpha(x, z) u(C_i^\alpha)^2 \right|$$

which by substitution of (18):

$$\leq \sum_{i=1}^n |C_i^\alpha(x)| |C_i^\alpha(z)| |w_\alpha(x, z)| \left( |C_i^\alpha(1)| \sqrt{n+1} \right)^{-2}$$

the maximum value of  $|C_i^\alpha(\cdot)| = |C_i^\alpha(1)|$ , thus:

$$\leq \sum_{i=1}^n |C_i^\alpha(1)|^2 |w_\alpha(x, z)| \left( |C_i^\alpha(1)| \sqrt{n+1} \right)^{-2} \quad (24)$$

$$\leq \sum_{i=1}^n \frac{|w_\alpha(x, z)|}{n+1}$$

the maximum of  $|w_\alpha(x_j, z_j)| = 1$ , therefore:

$$\leq \sum_{i=1}^n \frac{1}{n+1} < 1$$

Since the absolute value of the summation is always strictly less than 1, the summation can never be equal to  $-1$ . Further, for the

extreme values  $(x, z) \in \{(1, 1), (1, -1), (-1, 1), (-1, -1)\}$

$$\begin{aligned} & \sum_{i=0}^n C_i^\alpha(x_j) C_i^\alpha(z_j) w_\alpha(x_j, z_j) u(C_i^\alpha)^2 \\ &= \sum_{i=0}^n (\pm 1) w_\alpha(x_j, z_j) = \sum_{i=0}^n \pm \varepsilon \leq (n+1)\varepsilon \end{aligned} \quad (25)$$

which can be guaranteed to be  $\leq 1$  by requiring that  $\varepsilon \leq 1/(n+1)$ . ■

### 3.2.3. The Gegenbauer kernel prevents the adverse explosion effect

The explosion effect occurs when the product of polynomials  $C_n^\alpha(x)C_n^\alpha(z)$  powered by the problem dimension  $d$  produces large values that can lead to numerical issues. To prove that this effect is prevented, it must be proved that  $K_{Geg}(\mathbf{x}, \mathbf{z})$  is upper-bounded.

**Proof.** for any  $-1 \leq x \leq 1$ ,  $-1 \leq z \leq 1$ ,  $n \geq 0$ ,  $d \geq 1$  and  $\alpha > -0.5$ :

$$K_{Geg}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n C_i^\alpha(x_j) C_i^\alpha(z_j) w_\alpha(x_j, z_j) \left( |C_i^\alpha(1)| \sqrt{n+1} \right)^{-2}$$

following similar steps to those in (24):

$$K_{Geg}(\mathbf{x}, \mathbf{z}) \leq \prod_{j=1}^d \left( \sum_{i=0}^n \frac{1}{n+1} \right) \leq \prod_{j=1}^d \left( \frac{n+1}{n+1} \right) \leq 1 \quad (26)$$

Thus, it has been proved that the Gegenbauer kernel,  $K_{Geg}(\mathbf{x}, \mathbf{z}) \leq 1$ . ■

## 4. Experimental methodology

In this section, our aim is to compare the Gegenbauer kernel against classical kernels (Eqs. (27) to (29)) and previously described orthogonal polynomial kernels on the fairest ground possible, by means of individual hyper-parameter optimization and evaluation on a variety of datasets. The Hermite kernel formulation in Eq. (13) is very prone to the Explosion Effect, and was adapted under our formulation to the form in Eq. (30).

$$K_{Linear}(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z} \quad (27)$$

$$K_{RBF}(\mathbf{x}, \mathbf{z}) = e^{-\gamma \|\mathbf{x} - \mathbf{z}\|^2} \quad (28)$$

$$K_{Polynomial}(\mathbf{x}, \mathbf{z}) = (\mathbf{a}^T \mathbf{z} + b)^n \quad (29)$$

$$K_{S-Herm}(\mathbf{x}, \mathbf{z}) = \prod_{j=1}^d \sum_{i=0}^n He_i(x_j) He_i(z_j) \left( e^{-(x_j^2 + z_j^2)/2} \right) (2^{-2n}) \quad (30)$$

Fig. 5 shows the general process employed to evaluate seven state-of-the-art kernels (RBF, MCK, LINEAR, GCK, POLY, LEGEN and CHEB) against the two kernels following our formulation (GEGEN and s-HERM). Hereafter, these short labels will be used for the sake of presentation. The evaluation was carried out on 20 datasets for binary classification taken from the UCI Machine Learning,<sup>1</sup> LIBSVM,<sup>2</sup> and KEEL<sup>3</sup> repositories. These datasets, described in Table 2, were scaled to  $[-1, 1]$  in order to meet the requirements of orthogonal polynomial kernels and to avoid bias towards attributes of larger values.

Each of the nine kernels under evaluation has at least one hyper-parameter that needs to be tuned by the user to get the best

**Table 2**

Datasets description and its classification accuracy indexes.

Dataset <sup>(source)</sup>	#Instances	#Attributes	Reported accuracy (%) with RBF	[Ref.]
1 a1a <sup>2</sup>	1605	123	81.57	[32]
2 Australian <sup>2</sup>	690	14	81.87	[32]
3 Breast <sup>1</sup>	683	10	97.31	[33]
4 Diabetes <sup>1</sup>	768	8	77.73	[33]
5 Fourclass <sup>2</sup>	862	2	99.81	[33]
6 German <sup>2</sup>	1000	24	73.38	[32]
7 Glass <sup>2</sup>	214	9	92.06	[34]
8 Haberman <sup>1</sup>	306	3	73.55	[33]
9 Heart <sup>2</sup>	270	13	87.70	[15]
10 Ionosphere <sup>1</sup>	351	34	94.53	[15]
11 Liver <sup>1</sup>	345	7	72.45	[15]
12 Monks-1 <sup>1</sup>	124	6	95.01	[15]
13 Monks-2 <sup>1</sup>	169	6	76.58	[15]
14 Monks-3 <sup>1</sup>	122	6	89.37	[15]
15 plrx <sup>1</sup>	182	12	71.46	[35]
16 Sonar <sup>2</sup>	208	60	87.05	[15]
17 Splice <sup>2</sup>	1000	60	88.33	[36]
18 Vehicle <sup>2</sup>	1243	22	82.24	[37]
19 wdbc <sup>1</sup>	569	30	97.00	[32]
20 wpbc <sup>1</sup>	194	33	80.09	[15]

performance on a specific dataset. We performed an automatic hyper-parameter optimization by means of an efficient metaheuristic called BUMDA [31]. A hyper-parameter vector  $\mathbf{h} = (C, \gamma, a, n, \alpha)$  constitutes and individual of the population  $\mathbf{H}$  used by this metaheuristic. The specification of an initial population size (100 in this work) is the only parameter required by BUMDA at the beginning of the optimization. The ranges in which the best hyper-parameters were searched are the following:  $C \in (0, 2^5]$ ,  $n \in \{1, 2, \dots, 6\}$ ,  $\gamma \in (2^{-6}, 2^2]$ . The offset factor of the classical polynomial kernel was set to  $b = 0$  and its scale factor  $a$  was searched in the same range as  $\gamma$ . For the Gegenbauer kernel, the parameter  $\alpha$  was searched in  $(-0.5, 1.5]$ .

The quality of each kernel using a particular hyper-parameter setting was measured by the classification accuracy obtained by training an SVM with the LIBSVM solver [38] and evaluating it by 10-fold cross-validation. The same data folds and initial population  $\mathbf{H}^0$  of hyper-parameters were used by all the kernels to make the comparison between them as fair as possible. At the end of 15 generations, i.e. after evaluating around 1,500 different SVMs for each kernel, the optimization process ends by providing the best SVM kernel with its corresponding individual  $\mathbf{h}^*$  and performance measures (accuracy, PSV and training time). All the performance measures and all the evaluated SVM kernels are stored for subsequent analysis. To increase the confidence in our results, the optimization process was repeated 35 times for each dataset (35 experimental trials). Each experimental trial utilizes a different 10-fold data partition and a different  $\mathbf{H}^0$ . The average performance measures of these 35 optimal kernels are finally reported and were used for statistical tests. The Friedman, aligned Friedman and Quade tests were employed in order to draw stronger conclusions from the statistical point of view. The difference between these three methods is in the way of ranking each algorithm; a complete description of these tests is found in [39,40].

A suggested implementation of the Gegenbauer kernel functions for SVMs is provided in Algorithm 1. The implementation of the Pochhammer operator used in (18) is provided in Algorithm 2. Regarding the determination of parameters  $\alpha$  and  $n$ , any  $\alpha > -0.5$ , and  $n \in \mathbb{Z}^+$  can be used. However, in practice some extra considerations have to be taken into account. The orthogonal polynomial kernels are implemented by recurrence equations and the time required for their evaluation increases exponentially with the degree. In addition, in our experiments and in previous works, setting  $n \leq 6$  has sufficed to obtain good performance in a large number of prob-

<sup>1</sup> archive.ics.uci.edu/ml/datasets.html.

<sup>2</sup> www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html.

<sup>3</sup> sci2s.ugr.es/keel/dataset.php?cod=121.

Experimental Methodology	
INPUT: User-defined parameters (hyper parameter ranges, number of trials, etc.)	
OUTPUT: Performance Information (Accuracy and PSV) to carry out statistical analysis.	
1	FOR EACH experimental trial $\{i\}$ DO // $i : 1$ to 35 in this work
2	$\mathbf{H}_i^0 \leftarrow$ Generate a new initial population of hyper-parameters
3	FOR EACH $d \in \text{Datasets}$ DO
4	Partition $\text{Dataset}_d$ into $\text{DataFolds} = \{\text{Training}_{j,d}, \text{Test}_{j,d}\}_{j=1}^{10}$ //for 10-fold cross validation
5	FOR EACH $k \in \text{Kernels}$ DO
6	$\mathbf{H}_i^* \leftarrow \text{BUMDA}(\mathbf{H}_i^0, \text{Kernel}_k, \text{LIBSVM}, \text{DataFolds})$ // hyper parameter optimization
7	$[\text{Accuracy}_{i,k}, \text{PSV}_{i,k}] \leftarrow \text{SVM}(\text{Kernel}_k, \mathbf{H}_i^*)$ // store performance indexes
8	END
9	END
10	END

**Fig. 5.** Pseudocode of Experimental Methodology. Refer to the main text for explanation of this figure.

**Algorithm 1** The Gegenbauer kernel functions  $K_{\text{Geg}}(\mathbf{x}, \mathbf{z}, n, \alpha)$ .

INPUTS: $\mathbf{x}, \mathbf{z} \in [-1, 1]^d$ , $\alpha > -0.5$ , and $n \in \mathbb{Z}^+$	
OUTPUTS: $K_{\text{Geg}}(\mathbf{x}, \mathbf{z}, n, \alpha) \in \mathbb{R}$	
1	DEFINE: $\text{mult} \leftarrow 1$ , $\text{xlen} \leftarrow \text{lengthOf}(\mathbf{x})$ , $\text{ylen} \leftarrow \text{lengthOf}(\mathbf{z})$ , $i \leftarrow 0$ and $j \leftarrow 0$ ;
2	IF ( $\alpha == 0$ ) RETURN $C(x_i, k, 0)$ ; Eq. (16)
3	WHILE ( $i < \text{xlen}$ AND $j < \text{ylen}$ )
4	IF ( $\text{indexOf}(x_i) == \text{indexOf}(z_j)$ )
5	DEFINE: $\text{sum} \leftarrow 1$ ;
6	FOR ( $k = 1 : n$ ):
7	$\text{sum} \leftarrow \text{sum} + C(x_i, k, \alpha) \times C(z_j, k, \alpha) \times w(x_i, z_j, \alpha) \times u(k, \alpha)^2$ ; Eqs. (15, 18, 19)
8	END
9	$\text{mult} \leftarrow \text{mult} \times \text{sum}$ ; $i \leftarrow i + 1$ ; $j \leftarrow j + 1$ ;
10	ELSE
11	IF ( $\text{indexOf}(x_i) > \text{indexOf}(z_j)$ ) $j \leftarrow j + 1$ ;
12	ELSE $i \leftarrow i + 1$ ;
13	END
14	END
15	END
16	RETURN $\text{mult}$ ;

$$w(x_i, z_j, \alpha) = \begin{cases} 1, & -0.5 < \alpha \leq 0.5 \\ ((1 - x^2)(1 - z^2))^{\alpha - 1/2} + (\varepsilon = 0.1), & \alpha > 0.5 \end{cases} \quad \text{From Eq. (19)}$$

$$u(k, \alpha) = \beta \times |C_k^{\alpha}(1)| = (1/\sqrt{n+1}) \times (2\alpha)_k / (1)_k \quad \text{From Eq. (18)}$$

**Algorithm 2** Pochhammer Operator  $(x)_n$ .

INPUTS: $x \in \mathbb{R}$ , $n \in \mathbb{Z}$ ;	
OUTPUT: $(x \in \mathbb{R})_{n \in \mathbb{Z}} = (x)(x+1) \dots (x+n-1)$	
1	IF ( $n == 0$ ) RETURN 1.0;
2	IF ( $n < 0$ OR $x == 0$ ) RETURN 0.0;
3	DEFINE $\text{aux} \leftarrow 1$ ;
4	FOR ( $i = 1 : n - 1$ ): $\text{aux} \leftarrow \text{aux} \times (x + i)$ ; END
5	RETURN $\text{aux}$ ;

lems. Because of these reasons, we recommend to first test small degrees when using orthogonal polynomial kernels as done in our methodology. With respect to  $\alpha$ , we have noticed a tendency towards more training time required with values of  $\alpha$  larger than 2 and without a gain in performance. Thus, we recommend employing an initial setting of  $-0.5 < \alpha \leq 1.5$  for the Gegenbauer kernel.

## 5. Experimental results and discussion

Tables 3 and 4 summarize the results of the 35 trials of the experiment described in the previous section. In addition to the mean and standard deviation of the 35 best models found for each kernel, these tables include a ranking among the nine kernels. Average ranks measure the general performance of kernels and they are also used later on to apply the statistical tests. The performance of kernels in terms of accuracy can be observed in Table 3. The Gegenbauer kernel performs better than the rest, with the RBF and s-Hermite ranked 2nd and 3rd, respectively.

The general performance of the kernels in terms of PSV, complementing the accuracy index, is presented in Table 4. The linear and polynomial kernels obtain better ranks than the RBF kernel. The kernels s-Hermite and GCK, are ranked 1st and 2nd in PSV,



**Table 3**

Statistics on kernel performance according to classification accuracy.

Kernel	RBF		MCK		LINEAR		GCK		POLY		LEGEND		CHEB		s-HERM		GEGEN	
Dataset	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D	Mean	St.D
a1a	83.24 6	0.7	80.97 7	0.5	83.55 2	0.3	83.48 5	0.3	83.50 3.5	0.3	80.26 8	0.3	75.46 9	0.2	83.50 3.5	0.2	<b>83.64</b> 1	0.2
Australian	86.06 4	0.4	84.46 7	2.6	85.49 6	0.1	86.00 5	0.4	86.10 3	0.3	80.21 8	0.8	79.95 9	0.7	86.15 2	0.3	<b>86.51</b> 1	0.4
Breast	97.13 4	0.3	96.25 8	0.2	96.96 6	0.2	97.09 5	0.3	97.15 3	0.2	96.26 7	0.3	95.37 9	0.3	<b>97.37</b> 1	5.8	97.29 2	0.2
Diabetes	<b>78.04</b> 1	0.3	77.73 3	0.4	77.45 7	0.4	77.62 4	0.4	77.60 5	0.4	77.26 8	0.5	70.23 9	1.2	77.94 2	1.7	77.52 6	0.3
Fourclass	<b>100</b> 1.5	0.0	<b>100</b> 1.5	0.0	77.09 9	0.7	99.24 6	0.1	81.47 7	0.1	99.70 5	0.1	99.95 3	0.1	81.16 8	2.1	99.87 4	0.1
German	75.36 5	2.3	71.83 7	0.7	76.87 3	0.4	74.83 6	2.6	76.86 4	0.5	71.33 8	0.5	70.15 9	0.2	<b>77.20</b> 1	0.6	76.92 2	0.5
Glass 2	<b>92.65</b> 1	0.4	92.48 2	0.4	92.14 7.5	0.0	92.35 3	0.3	92.33 4	0.3	92.20 5	0.3	89.49 9	0.9	92.14 7.5	0.4	92.19 6	0.3
Haberman	74.81 3	0.6	75.02 2	0.7	73.42 8	0.3	74.43 6	0.5	74.21 7	0.7	74.62 4	0.4	72.44 9	0.9	74.44 5	0.7	<b>75.09</b> 1	0.6
Heart	84.39 3	0.5	81.43 7	2.0	84.21 4	0.9	82.02 6	1.4	83.94 5	0.6	77.56 8	0.8	72.72 9	0.9	85.26 2	2.3	<b>85.98</b> 1	0.5
Ionosphere	<b>95.88</b> 1	0.3	95.04 2	0.4	89.38 8	0.7	91.89 7	0.6	92.24 6	0.6	92.82 5	0.6	69.22 9	8.2	94.12 3	4.0	93.37 4	0.4
Liver	74.25 3	0.6	<b>74.85</b> 1	0.7	69.62 9	0.9	74.71 2	0.6	72.80 7	0.6	73.08 6	0.7	71.47 8	0.8	73.89 5	2.0	74.13 4	0.7
Monks-1	82.56 4	1.6	87.75 2	2.3	69.20 8	2.8	74.69 6	2.4	72.23 7	1.9	76.33 5	1.8	68.99 9	1.3	85.75 3	2.0	<b>97.52</b> 1	0.9
Monks-2	81.95 4	1.7	<b>94.55</b> 1	3.5	62.15 9	0.0	77.75 5	2.0	71.77 8	1.6	85.03 3	1.6	73.19 7	1.8	73.58 6	6.5	85.64 2	1.5
Monks-3	91.88 3	0.8	91.47 4	1.1	81.44 7	1.6	88.31 5	1.1	80.12 8	2.0	84.03 6	2.1	73.42 9	1.9	93.42 2	7.1	<b>93.46</b> 1	0.1
plrx	73.40 2	0.8	<b>73.67</b> 1	0.9	71.46 9	0.0	72.43 4	0.3	71.97 7	0.5	71.53 8	0.3	72.85 3	0.9	72.42 5	0.5	72.31 6	0.4
Sonar	90.53 2	1.0	89.92 3	1.0	80.08 7	1.1	86.68 6	1.0	87.29 5	1.2	73.55 8	0.9	70.21 9	0.8	89.91 4	1.6	<b>92.25</b> 1	1.0
Splice	87.96 4	5.0	89.08 3	0.5	80.64 7	0.8	87.22 5	0.5	82.94 6	0.5	56.92 8	0.4	51.83 9	1.3	92.26 2	0.4	<b>93.62</b> 1	0.3
Vehicle	84.56 6	3.6	84.69 5	1.5	82.73 7	0.2	<b>85.51</b> 1	0.3	85.07 3	0.3	77.64 8	0.5	75.09 9	0.6	85.05 4	0.7	85.24 2	0.4
wdbc	98.23 2	0.2	97.24 7	0.4	97.87 6	0.2	98.09 4	0.2	97.92 5	0.2	95.45 9	0.4	95.72 8	0.4	<b>98.35</b> 1	0.2	98.15 3	0.3
wdbc	82.13 2	1.5	78.32 7	1.1	81.51 5	1.0	81.40 6	1.2	81.80 4	1.2	76.40 9	0.2	76.54 8	0.4	82.11 3	1.6	<b>82.56</b> 1	1.0
Avg. Rank	3.075		4.025		6.725		4.85		5.375		6.8		8.15		3.5		<b>2.5</b>	

but the s-Hermite kernel is ranked 3rd in accuracy while GCK is 5th.

The training times required for each combination of a kernel and a dataset are reported in Table 5. These values can be analyzed by groups based on accuracy: the top-3 kernels are Gegenbauer, RBF and s-Hermite; the second group, with medium performance includes the MCK, GCK and Polynomial kernels; the bottom-3 kernels are Linear, Legendre and Chebyshev. It can be observed that the top-3 kernels, with an average total time required of 44.5 (s), outperform the second group which requires an average of 236.2 (s) and also the bottom-3 group with 174.3 (s).

### 5.1. Comparison between all kernels

Results on accuracy and PSV reported in Tables 3 and 4, respectively, were obtained by means of an optimization of the parameters corresponding to each kernel and the hyper-parameter C, which represents the penalty factor in SVMs. In general, it was found that the kernel parameters have a stronger influence on performance of the SVMs than C. In every case, all measures were taken to ensure that the evaluation was as independent as possible from the choice of the parameter values, so that the true potential of each kernel could be appreciated.

The Gegenbauer, RBF and s-Hermite kernels emerge as the best three out of the nine considered in terms of accuracy. To verify this

observation and to get a closer view of their differences in performance, the statistical tests of Friedman, Aligned Friedman and Quade were applied. The results of these tests, performed under a significance level of 0.05, are reported in Tables 6–8.

The statistical tests performed look for a significant difference between at least two kernels. Any  $p$ -value  $< 0.05$  indicates that there exist at least two kernels with different performance. The results in Table 6 show that, according to the Friedman, Aligned Friedman and Quade tests, the performance between some kernels differ significantly with respect to the accuracy index ( $p$ -values 5.29E-11, 0.023 and 1.8E-13). According to the Aligned Friedman test, there are also differences among the performance of the kernels with respect to the PSV index ( $p$ -value 0.02). After *post hoc* tests of Bonferroni-Dunn, Holm, Holland and Finner [39] using Gegenbauer and GCK as the control methods under accuracy and PSV indexes, respectively, it was found that Gegenbauer is statistically superior in accuracy to all kernels except for RBF and s-Hermite and that the GCK is significantly better in PSV than the Legendre kernel.

In order to obtain a deeper understanding of the relative behavior of the proposed kernel, we performed separate statistical tests comparing the Gegenbauer kernel against the classical kernels and against the other orthogonal polynomial kernels. The results of these comparisons are discussed below.

**Table 4**  
Statistics on kernels performance according to proportion of support vectors.

Kernel	RBF		MCK		LINEAR		GCK		POLY		LEGEN		CHEBY		s-HERM		GEGEN	
Dataset	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv	Mean	St.Dv
a1a	42.9 7	7.21	61.7 8	10.1	36.6 2	1.35	39.41 4	1.44	36.8 3	1.5	76.3 9	0.09	<b>3.10</b> 1	2.27	41.0 6	0.79	40.2 5	1.25
Australian	38.4 6	9.16	39.6 7	11.6	<b>31.3</b> 1	0.55	35.99 3	1.95	37.0 5	14.2	55.0 9	12.8	43.9 8	0.17	35.3 2	7.92	36.3 4	8.70
Breast	13.3 3	8.07	36.7 9	9.55	<b>8.37</b> 1	1.32	9.38 2	2.36	14.0 4	11.7	26.6 7	4.54	28.2 8	3.70	14.26 5	5.80	14.4 6	5.44
Diabetes	55.4 6	3.17	56.4 9	4.87	54.10 4	2.88	55.63 7	3.94	55.64 8	4.35	51.7 2	1.96	<b>48.63</b> 1	2.44	53.76 3	1.75	55.1 5	4.23
Fourclass	6.30 3.5	2.42	<b>4.17</b> 1	0.92	55.17 9	5.84	9.01 5	1.78	38.4 7	0.85	6.30 3.5	1.06	4.88 2	0.51	51.70 8	2.17	10.5 6	0.90
German	63.9 6	12.9	84.4 7	8.7	<b>53.4</b> 1	0.55	61.08 5	0.56	53.8 2	2.38	89.8 8	1.38	99.0 9	0.62	57.2 4	1.11	56.4 3	1.28
Glass2	29.6 8	7.69	25.1 6	7.62	27.10 7	0.67	20.61 2	1.59	24.4 5	3.45	<b>16.64</b> 1	2.16	39.7 9	9.77	21.43 3	0.42	22.6 4	0.58
Haberman	54.4 6	0.59	53.8 2	0.45	54.32 5	0.27	54.78 7	0.42	54.8 8	0.99	<b>52.51</b> 1	0.47	55.4 9	0.62	54.1 4	0.77	53.9 3	1.41
Heart	44.7 4	6.25	45.2 5	8.04	39.52 2	4.65	46.92 6	6.30	<b>37.76</b> 1	2.45	66.1 9	2.19	62.9 8	0.18	43.01 3	2.34	47.5 7	3.78
Ionosphere	41.6 7	12.0	51.5 9	9.60	27.11 4	5.85	26.33 3	1.78	28.4 5	7.56	41.8 8	0.20	<b>21.43</b> 1	3.98	24.92 2	4.06	32.7 6	2.91
Liver	71.3 7	3.95	66.0 3	4.06	75.79 9	2.07	68.31 4	4.18	68.5 5	3.73	65.8 2	4.46	<b>62.83</b> 1	3.61	73.46 8	2.08	69.2 6	5.25
Monks-1	53.3 3	2.81	<b>33.7</b> 1	6.50	74.21 9	10.1	56.16 5	13.15	62.2 6	7.54	66.8 7	3.47	70.8 8	1.62	53.69 4	2.02	47.2 2	3.82
Monks-2	73.8 8	11.2	<b>34.8</b> 1	9.97	91.15 9	0.38	57.70 5	1.30	55.3 4	5.22	42.1 3	1.02	70.9 7	2.30	60.21 6	6.52	40.9 2	2.04
Monks-3	47.5 3	5.49	<b>38.5</b> 1	9.52	52.7 4.5	4.81	58.21 8	6.97	54.8 7	9.10	52.7 4.5	9.96	62.7 9	3.24	44.56 2	7.16	53.9 6	8.78
plrx	94.9 7	4.28	96.2 8	0.76	70.28 3	0.94	70.55 4	0.86	72.2 5	1.89	<b>59.14</b> 1	8.23	96.6 9	0.07	70.24 2	0.56	72.6 6	3.22
Sonar	67.7 7	8.35	60.5 5	5.95	60.07 4	10.7	<b>43.92</b> 1	2.80	57.4 2	3.35	69.9 8	0.25	77.7 9	0.21	58.29 3	1.67	64.0 6	3.20
Splice	69.2 5	8.73	79.82 7	8.67	43.93 2	7.19	56.40 3	2.71	92.7 8	0.14	97.7 9	0.02	<b>1.99</b> 1	1.72	71.53 6	0.45	61.8 4	3.35
Vehicle	42.0 8	8.18	39.11 4	1.90	41.55 7	0.24	39.10 3	0.70	<b>33.83</b> 1	1.38	60.9 9	10.5	38.7 2	0.15	40.18 6	0.74	39.4 5	0.69
wdbc	12.4 9	1.76	11.2 6	4.49	11.5 8	0.22	11.41 7	1.67	10.64 4	2.86	10.9 5	3.57	10.4 2	0.08	<b>9.42</b> 1	1.31	10.62 3	2.01
wdbc	57.9 7	9.57	81.37 8	14.0	<b>45.12</b> 1	1.84	50.85 3	1.10	45.6 2	6.06	51.2 4	10.6	96.7 9	5.34	55.45 5	1.64	55.7 6	2.44
Avg. Rank	6.025		5.35		4.625		4.35		4.6		5.5		5.65		<b>4.15</b>		4.75	

**Table 5**  
Training times (in seconds) of each kernel using its best parameters, Mean  $\pm$  Std. Dev.

Dataset	RBF	MCK	LINEAR	GCK	POLY	LEGEN	CHEB	s-HERM	GEGEN
a1a	0.22 $\pm$ 0.1	38.91 $\pm$ 31.5	12.85 $\pm$ 5.1	6.82 $\pm$ 5.1	36.53 $\pm$ 25.3	1.13 $\pm$ 0.1	238.5 $\pm$ 28.9	7.33 $\pm$ 2.0	4.10 $\pm$ 3.9
Australian	0.04 $\pm$ 0.0	0.87 $\pm$ 3.6	3.57 $\pm$ 2.5	0.34 $\pm$ 0.3	0.26 $\pm$ 0.4	0.92 $\pm$ 0.5	0.41 $\pm$ 0.1	4.00 $\pm$ 2.4	2.88 $\pm$ 1.7
Breast	0.04 $\pm$ 0.1	6.48 $\pm$ 6.2	0.28 $\pm$ 0.3	0.17 $\pm$ 0.3	0.24 $\pm$ 0.3	0.32 $\pm$ 0.2	0.34 $\pm$ 0.3	0.40 $\pm$ 0.3	0.94 $\pm$ 1.6
Diabetes	0.26 $\pm$ 0.1	8.57 $\pm$ 33.1	0.28 $\pm$ 0.3	0.68 $\pm$ 1.0	5.75 $\pm$ 10.0	0.64 $\pm$ 0.3	2.65 $\pm$ 1.1	5.73 $\pm$ 3.6	5.21 $\pm$ 6.2
Fourclass	0.09 $\pm$ 0.2	1.94 $\pm$ 1.2	0.75 $\pm$ 0.6	3.05 $\pm$ 0.8	239.9 $\pm$ 192.0	0.66 $\pm$ 0.3	0.46 $\pm$ 0.4	4.38 $\pm$ 1.0	1.45 $\pm$ 0.4
German	0.28 $\pm$ 0.0	50.84 $\pm$ 44.9	7.85 $\pm$ 3.6	0.83 $\pm$ 0.7	4.14 $\pm$ 4.7	0.82 $\pm$ 0.1	1.40 $\pm$ 0.1	4.96 $\pm$ 3.9	4.90 $\pm$ 4.3
Glass2	0.02 $\pm$ 0.1	0.77 $\pm$ 1.5	0.003 $\pm$ 0.0	1.07 $\pm$ 1.3	0.34 $\pm$ 0.8	0.02 $\pm$ 0.0	0.36 $\pm$ 0.4	0.17 $\pm$ 0.2	0.22 $\pm$ 0.2
Haberman	0.09 $\pm$ 0.2	2.62 $\pm$ 14.5	0.07 $\pm$ 0.2	0.30 $\pm$ 0.4	40.19 $\pm$ 43.1	1.01 $\pm$ 1.0	0.38 $\pm$ 0.3	0.86 $\pm$ 0.5	1.08 $\pm$ 0.9
Heart	0.05 $\pm$ 0.1	1.81 $\pm$ 1.7	0.62 $\pm$ 0.4	1.27 $\pm$ 1.6	3.90 $\pm$ 2.2	0.31 $\pm$ 0.2	0.12 $\pm$ 0.2	1.30 $\pm$ 1.3	1.13 $\pm$ 0.6
Ionosphere	0.05 $\pm$ 0.1	0.57 $\pm$ 0.8	0.62 $\pm$ 0.7	0.22 $\pm$ 0.3	0.13 $\pm$ 0.2	0.71 $\pm$ 0.4	51.30 $\pm$ 22.3	2.55 $\pm$ 1.4	5.10 $\pm$ 3.9
Liver	0.09 $\pm$ 0.1	1.52 $\pm$ 10.0	0.29 $\pm$ 0.4	0.32 $\pm$ 0.6	3.74 $\pm$ 2.8	1.88 $\pm$ 1.1	0.57 $\pm$ 0.5	1.18 $\pm$ 1.9	0.78 $\pm$ 0.7
Monks-1	0.005 $\pm$ 0.0	0.79 $\pm$ 1.2	0.003 $\pm$ 0.0	0.12 $\pm$ 0.1	0.03 $\pm$ 0.1	0.01 $\pm$ 0.0	0.02 $\pm$ 0.0	0.003 $\pm$ 0.0	0.11 $\pm$ 0.0
Monks-2	0.11 $\pm$ 0.3	0.39 $\pm$ 1.0	0.04 $\pm$ 0.1	0.06 $\pm$ 0.1	0.92 $\pm$ 0.9	0.23 $\pm$ 0.3	0.02 $\pm$ 0.0	0.08 $\pm$ 0.2	0.26 $\pm$ 0.1
Monks-3	0.05 $\pm$ 0.2	0.21 $\pm$ 0.2	0.004 $\pm$ 0.0	0.02 $\pm$ 0.0	0.02 $\pm$ 0.1	0.05 $\pm$ 0.1	0.06 $\pm$ 0.1	0.14 $\pm$ 0.3	0.30 $\pm$ 0.2
plrx	0.15 $\pm$ 0.3	1.09 $\pm$ 1.1	1.22 $\pm$ 1.0	0.95 $\pm$ 0.9	0.05 $\pm$ 0.2	0.03 $\pm$ 0.1	0.57 $\pm$ 0.5	0.63 $\pm$ 0.5	0.72 $\pm$ 0.5
Sonar	0.09 $\pm$ 0.1	4.27 $\pm$ 3.9	0.10 $\pm$ 0.2	2.48 $\pm$ 2.7	0.11 $\pm$ 0.2	0.58 $\pm$ 0.5	0.60 $\pm$ 0.5	0.25 $\pm$ 0.3	3.63 $\pm$ 1.7
Splice	0.37 $\pm$ 0.1	13.89 $\pm$ 26.5	0.47 $\pm$ 0.3	9.78 $\pm$ 5.5	179.1 $\pm$ 131.4	3.44 $\pm$ 1.8	179.5 $\pm$ 55.9	17.55 $\pm$ 7.7	5.29 $\pm$ 3.8
Vehicle	0.34 $\pm$ 0.2	12.47 $\pm$ 85.6	0.83 $\pm$ 1.0	3.23 $\pm$ 4.2	4.82 $\pm$ 4.3	2.11 $\pm$ 1.1	0.01 $\pm$ 0.0	17.13 $\pm$ 8.1	21.42 $\pm$ 15.4
wdbc	0.02 $\pm$ 0.0	4.73 $\pm$ 4.1	0.03 $\pm$ 0.0	0.47 $\pm$ 0.3	0.03 $\pm$ 0.0	0.19 $\pm$ 0.1	0.08 $\pm$ 0.1	0.56 $\pm$ 0.6	1.25 $\pm$ 0.9
wdbc	0.06 $\pm$ 0.1	2.04 $\pm$ 2.9	0.15 $\pm$ 0.2	0.22 $\pm$ 0.2	1.11 $\pm$ 0.7	0.10 $\pm$ 0.2	0.21 $\pm$ 0.2	0.47 $\pm$ 0.5	0.68 $\pm$ 0.1
Sum of means	2.41	154.79	30.05	32.41	521.35	15.18	477.55	69.70	61.45

**Table 6**  
Average ranks and *p*-values of statistical tests on all kernels.

Algorithm	Friedman		Aligned Friedman		Quade	
	Accuracy	PSV	Accuracy	PSV	Accuracy	PSV
RBF	3.075	6.025	55.775	104.425	3.176	5.876
MCK	4.025	5.350	73.775	103.200	3.619	5.424
LINEAR	6.725	4.625	116.475	84.075	6.855	4.283
GCK	4.850	4.350	82.100	<b>69.000</b>	5.110	<b>4.014</b>
POLYNOMIAL	5.375	4.600	94.175	81.650	5.698	4.338
LEGENDRE	6.800	5.500	123.000	107.950	6.548	6.117
CHEBYSHEV	8.150	5.650	148.950	107.550	8.224	5.605
s-HERMITE	3.500	<b>4.150</b>	69.800	80.150	3.576	4.600
GEGBAUER	<b>2.500</b>	4.750	<b>50.450</b>	76.500	<b>2.195</b>	4.743
<i>p</i> -value	<b>8.16E-11</b>	0.35	<b>0.023</b>	<b>0.02</b>	<b>2.87E-14</b>	0.29

**Table 7**  
Average ranks and *p*-values of tests on classical kernels vs Gegenbauer.

Algorithm	Friedman		Aligned Friedman		Quade	
	Accuracy	PSV	Accuracy	PSV	Accuracy	PSV
RBF	2.000	3.200	26.45	53.00	1.880	3.123
LINEAR	3.700	<b>2.100</b>	62.15	37.10	3.761	<b>2.262</b>
POLYNOMIAL	2.850	2.400	52.35	<b>35.80</b>	3.000	2.300
GEGBAUER	<b>1.450</b>	2.300	<b>21.05</b>	36.10	<b>1.357</b>	2.314
<i>p</i> -value	<b>1.23E-07</b>	<b>0.038</b>	<b>0.0013</b>	<b>9.23E-04</b>	<b>9.30E-15</b>	<b>0.001</b>

**Table 8**  
Average ranks and *p*-values of statistical tests on OP kernels vs Gegenbauer.

Algorithm	Friedman		Aligned Friedman		Quade	
	Accuracy	PSV	Accuracy	PSV	Accuracy	PSV
MCK	2.650	3.650	47.100	67.350	2.500	3.633
GCK	3.400	3.100	53.900	<b>47.450</b>	3.671	<b>2.738</b>
LEGENDRE	4.650	3.800	81.600	71.600	4.586	4.043
CHEBYSHEV	5.500	3.950	98.950	72.950	5.643	4.029
s-HERMITE	2.750	<b>3.050</b>	46.200	52.500	2.786	3.210
GEGBAUER	<b>2.050</b>	3.450	<b>35.250</b>	51.150	<b>1.814</b>	3.348
<i>p</i> -value	<b>1.53E-09</b>	0.566	<b>0.005</b>	<b>0.004</b>	<b>1.60E-13</b>	0.028

### 5.2. Comparison of Gegenbauer versus classical kernels

The Gegenbauer kernel was compared against the linear, polynomial and RBF kernels. The results of this comparison are presented in Table 7. In terms of Accuracy, it can be observed that the Gegenbauer kernel achieved better performance than the linear and polynomial kernels. Taking the Gegenbauer kernel as a control method and carrying out *post hoc* tests it was confirmed that Gegenbauer is significantly better than the linear and polynomial kernels but not significantly different than the RBF kernel. In terms of PSV, Gegenbauer is not different than the linear and polynomial kernels, but it is better, with statistical significance, than the RBF kernel.

### 5.3. Comparison of Gegenbauer versus orthogonal polynomial kernels

Having compared the Gegenbauer kernel against classical kernels, now we compare Gegenbauer against the orthogonal polynomial kernels. The results, reported in Table 8, show that with respect to accuracy, Gegenbauer was ranked first in all the tests. *Post hoc* tests revealed that the Gegenbauer kernel is better than all orthogonal polynomial kernels except for our modified kernel s-Hermite and MCK. In terms of PSV, the *post hoc* tests show that Gegenbauer is better, with statistical significance, than the Chebyshev kernel only.

As a general conclusion, the Gegenbauer and the s-Hermite kernels are generally better in classification accuracy than the rest of

the orthogonal polynomial kernels and maintain a low PSV. The RBF obtains the second best rank overall according to the accuracy index, but associated with the worst rank overall according to the PSV index. In accordance with results of previous studies [13,14], radial kernels in general require a larger amount of support vectors than the orthogonal polynomial kernels. The theory indicates that this may be detrimental to their generalization capabilities.

### 5.4. Hyper-parameter optimization results

The best values for the different parameters of the kernels discussed in this work are reported in Tables 9–11 and are the result of the optimization performed as part of the experimental methodology described above. The main focus in this work is not hyper-parameter optimization, however, Tables 9–11 contain important information that can be useful as a reference for those working in that problem, as well as to ensure the reproducibility of our study. The best parameter values were computed based on the results of the best 25% of the SVMs population, sorted by accuracy.

As can be appreciated, the values reported in these tables present a large variability and because of this, it becomes difficult to describe them in brief. Nevertheless, an important observation to be made is that a large standard deviation (relative to the range in which each parameter was optimized), indicates that the parameter has a weak influence in the performance of the kernel on the corresponding dataset. In the case of the polynomial degree in Table 11, a large proportion of occurrence of a particular degree indicates that the degree is the best choice of this parameter for the dataset, while a small proportion indicates that different degrees can produce very similar performance and the reported value corresponds to the majority of cases.

Table 9 contains the best values of the parameter *C* for each combination of a kernel and a dataset. This parameter is required by the SVM classifiers independently of the type of kernel function. As can be observed, for a significant number of cases, the standard deviation obtained is relatively large when compared with the range on which this parameter was optimized. This observation indicates that *C* has a lower impact on performance than other parameters, as will be shown next.

Table 10 shows the best value found for the intrinsic parameters of kernels, these are:  $\gamma$  (RBF and MCK),  $\alpha$  (polynomial) and  $\alpha$  (Gegenbauer). In this table, the reported standard deviations are relatively small, which unlike the case of the *C* parameter, indicates that the intrinsic kernel parameters possess a stronger influence on the performance of SVM classifiers.

Table 11 contains the best value of the polynomial degree  $n$  for the polynomial kernels (GCK, Legendre, Chebyshev, s-Hermite, MCK, Polynomial and Gegenbauer); the values appearing within parentheses represent the highest percentage of occurrence of a particular degree within the top 25% of the evaluated SVMs, sorted by accuracy and without duplicates. A relatively small value, such as (0.20) for MCK on dataset ‘a1a’, indicates large diversity in the parameters of the best solutions (i.e. other degrees are close in performance to the optimum); in contrast, a value close to 1, such as (0.96) for MCK on dataset ‘Australian’, indicates dominance of the corresponding degree within the best solutions found. Please notice that the results in Table 11 should not be interpreted in isolation from other results. For instance, the fact that the best degree found for the Polynomial kernel on dataset ‘Monks-3’ is  $n = 1$  with a high proportion of occurrence (0.87), does not necessarily mean that said dataset is easily solved by an essentially linear kernel. This becomes apparent by looking at the classification accuracy of the Polynomial kernel (cf. Table 3), according to which the kernel is ranked second from last, for that same dataset.

Interestingly, in Table 11 the smaller degrees (1, 2 and 3) appear most often than the larger degrees (4, 5 and 6), both globally

**Table 9**Best values of the C parameter, Mean  $\pm$  Std. Dev.

Dataset	RBF	MCK	LINEAR	GCK	POLY	LEGEN	CHEBY	s-HERM	GEGEN
a1a	10.45 $\pm$ 3.7	15.18 $\pm$ 7.72	22.34 $\pm$ 9.7	8.66 $\pm$ 8.07	18.65 $\pm$ 7.5	15.12 $\pm$ 8.9	18.77 $\pm$ 7.9	12.21 $\pm$ 6.7	24.46 $\pm$ 3.6
Australian	14.25 $\pm$ 6.2	5.06 $\pm$ 5.03	10.88 $\pm$ 9.2	8.85 $\pm$ 8.15	14.44 $\pm$ 7.6	2.22 $\pm$ 3.6	18.84 $\pm$ 9.0	18.70 $\pm$ 6.8	24.28 $\pm$ 4.8
Breast	19.34 $\pm$ 1.1	15.27 $\pm$ 9.21	16.31 $\pm$ 12.2	5.19 $\pm$ 5.24	17.38 $\pm$ 8.1	16.09 $\pm$ 9.6	17.42 $\pm$ 9.3	7.44 $\pm$ 4.0	16.25 $\pm$ 8.3
Diabetes	11.83 $\pm$ 0.2	7.41 $\pm$ 6.43	6.70 $\pm$ 0.6	10.19 $\pm$ 8.26	13.92 $\pm$ 7.2	5.43 $\pm$ 3.1	3.91 $\pm$ 1.2	16.43 $\pm$ 7.5	18.75 $\pm$ 8.4
Fourclass	30.42 $\pm$ 1.5	15.72 $\pm$ 9.86	11.77 $\pm$ 4.3	18.98 $\pm$ 7.67	21.63 $\pm$ 7.7	31.58 $\pm$ 1.0	28.85 $\pm$ 1.9	25.20 $\pm$ 4.4	31.52 $\pm$ 0.7
German	13.12 $\pm$ 2.0	15.93 $\pm$ 8.84	18.67 $\pm$ 9.2	0.75 $\pm$ 0.69	13.46 $\pm$ 8.3	15.761 $\pm$ 0.0	13.51 $\pm$ 8.5	3.18 $\pm$ 3.1	16.89 $\pm$ 7.1
Glass2	23.49 $\pm$ 1.4	11.99 $\pm$ 9.73	0.29 $\pm$ 0.3	10.03 $\pm$ 9.38	15.46 $\pm$ 9.1	17.98 $\pm$ 8.2	16.65 $\pm$ 11.4	16.65 $\pm$ 9.9	11.07 $\pm$ 11.0
Haberman	12.52 $\pm$ 1.1	9.38 $\pm$ 7.14	0.79 $\pm$ 0.4	10.72 $\pm$ 8.53	20.51 $\pm$ 6.3	15.56 $\pm$ 9.5	1.85 $\pm$ 0.5	25.26 $\pm$ 3.6	21.13 $\pm$ 5.3
Heart	9.64 $\pm$ 0.1	12.68 $\pm$ 10.01	9.35 $\pm$ 1.0	6.23 $\pm$ 9.59	19.18 $\pm$ 4.9	14.24 $\pm$ 9.3	13.42 $\pm$ 7.6	9.65 $\pm$ 2.4	11.58 $\pm$ 6.9
Ionosphere	18.30 $\pm$ 8.2	14.33 $\pm$ 7.97	30.04 $\pm$ 1.5	9.51 $\pm$ 8.27	10.57 $\pm$ 5.1	19.31 $\pm$ 6.3	12.43 $\pm$ 4.0	10.77 $\pm$ 4.3	5.32 $\pm$ 3.7
Liver	29.08 $\pm$ 3.0	10.81 $\pm$ 7.05	20.37 $\pm$ 0.1	14.64 $\pm$ 8.52	12.55 $\pm$ 5.0	15.22 $\pm$ 6.4	22.11 $\pm$ 3.9	28.53 $\pm$ 3.1	25.30 $\pm$ 4.8
Monks-1	16.86 $\pm$ 0.0	14.61 $\pm$ 6.89	1.00 $\pm$ 1.0	11.80 $\pm$ 11.29	6.24 $\pm$ 1.0	15.94 $\pm$ 8.1	18.20 $\pm$ 8.4	20.61 $\pm$ 0.5	20.96 $\pm$ 3.4
Monks-2	21.64 $\pm$ 2.4	15.96 $\pm$ 7.01	0.21 $\pm$ 0.1	12.92 $\pm$ 7.76	12.70 $\pm$ 6.6	5.44 $\pm$ 4.4	17.59 $\pm$ 9.1	29.84 $\pm$ 2.0	18.06 $\pm$ 5.8
Monks-3	29.17 $\pm$ 1.4	13.15 $\pm$ 8.20	1.37 $\pm$ 0.2	1.62 $\pm$ 2.14	6.13 $\pm$ 1.8	0.05 $\pm$ 0.0	10.55 $\pm$ 6.9	29.42 $\pm$ 1.8	27.04 $\pm$ 5.0
plrx	23.63 $\pm$ 7.9	15.45 $\pm$ 9.86	21.51 $\pm$ 9.0	10.94 $\pm$ 10.07	14.99 $\pm$ 7.8	13.72 $\pm$ 9.1	21.43 $\pm$ 11.2	23.13 $\pm$ 8.0	12.68 $\pm$ 9.1
Sonar	21.93 $\pm$ 2.4	14.16 $\pm$ 7.46	0.42 $\pm$ 0.0	14.18 $\pm$ 8.54	20.43 $\pm$ 7.9	16.96 $\pm$ 8.6	19.58 $\pm$ 9.8	20.44 $\pm$ 8.6	17.60 $\pm$ 9.3
Splice	9.77 $\pm$ 1.2	14.47 $\pm$ 6.73	0.88 $\pm$ 0.7	9.64 $\pm$ 9.76	17.08 $\pm$ 7.8	14.27 $\pm$ 7.6	13.75 $\pm$ 10.1	19.95 $\pm$ 9.8	19.03 $\pm$ 7.9
Vehicle	28.55 $\pm$ 1.8	9.69 $\pm$ 6.27	17.15 $\pm$ 4.4	14.31 $\pm$ 9.48	12.58 $\pm$ 7.2	27.45 $\pm$ 6.1	18.75 $\pm$ 8.7	22.24 $\pm$ 6.0	20.20 $\pm$ 7.5
wdbc	22.39 $\pm$ 3.6	15.01 $\pm$ 8.47	10.04 $\pm$ 2.5	8.01 $\pm$ 8.00	14.68 $\pm$ 9.1	16.95 $\pm$ 9.2	22.03 $\pm$ 5.3	24.64 $\pm$ 5.4	20.23 $\pm$ 8.4
wdbc	27.84 $\pm$ 6.2	14.76 $\pm$ 8.74	6.59 $\pm$ 2.2	13.65 $\pm$ 9.36	17.23 $\pm$ 3.1	14.12 $\pm$ 10.3	15.02 $\pm$ 8.7	10.07 $\pm$ 5.7	20.48 $\pm$ 3.6

**Table 10**Best value of intrinsic kernel parameter, Mean  $\pm$  Std. Dev.

Dataset	RBF $\gamma$	MCK $\gamma$	POLY $\alpha$	GEGEN $\alpha$
a1a	0.02 $\pm$ 0.0	0.23 $\pm$ 0.08	2.06 $\pm$ 1.0	0.96 $\pm$ 0.2
Australian	0.06 $\pm$ 0.0	0.31 $\pm$ 0.35	0.30 $\pm$ 0.2	1.20 $\pm$ 0.2
Breast	0.03 $\pm$ 0.0	1.94 $\pm$ 0.76	1.55 $\pm$ 1.1	0.50 $\pm$ 0.3
Diabetes	0.18 $\pm$ 0.0	0.34 $\pm$ 0.45	1.71 $\pm$ 1.0	0.61 $\pm$ 0.3
Fourclass	3.82 $\pm$ 0.1	1.12 $\pm$ 0.68	3.33 $\pm$ 0.3	-0.42 $\pm$ 0.1
German	0.02 $\pm$ 0.0	0.41 $\pm$ 0.24	2.01 $\pm$ 1.0	0.85 $\pm$ 0.4
Glass2	3.13 $\pm$ 0.4	1.77 $\pm$ 1.29	0.58 $\pm$ 0.9	0.38 $\pm$ 0.3
Haberman	0.71 $\pm$ 0.0	0.54 $\pm$ 0.67	2.97 $\pm$ 0.5	0.00 $\pm$ 0.2
Heart	0.02 $\pm$ 0.0	0.24 $\pm$ 0.21	2.55 $\pm$ 0.7	0.84 $\pm$ 0.2
Ionosphere	0.05 $\pm$ 0.0	0.70 $\pm$ 0.50	0.14 $\pm$ 0.0	1.36 $\pm$ 0.1
Liver	0.59 $\pm$ 0.2	0.39 $\pm$ 0.44	1.44 $\pm$ 0.2	0.23 $\pm$ 0.1
Monks-1	0.08 $\pm$ 0.0	0.05 $\pm$ 0.05	0.34 $\pm$ 0.0	0.38 $\pm$ 0.0
Monks-2	0.43 $\pm$ 0.0	0.15 $\pm$ 0.27	2.39 $\pm$ 1.0	-0.44 $\pm$ 0.1
Monks-3	0.08 $\pm$ 0.0	0.17 $\pm$ 0.22	0.36 $\pm$ 0.2	1.11 $\pm$ 0.4
plrx	0.05 $\pm$ 0.0	3.20 $\pm$ 0.72	0.02 $\pm$ 0.0	0.51 $\pm$ 0.0
Sonar	0.04 $\pm$ 0.0	0.11 $\pm$ 0.07	1.69 $\pm$ 1.1	0.74 $\pm$ 0.1
Splice	0.02 $\pm$ 0.0	0.04 $\pm$ 0.02	2.25 $\pm$ 1.2	0.74 $\pm$ 0.1
Vehicle	1.27 $\pm$ 0.0	0.16 $\pm$ 0.13	1.32 $\pm$ 0.3	0.76 $\pm$ 0.4
wdbc	0.12 $\pm$ 0.0	0.41 $\pm$ 0.26	0.62 $\pm$ 0.5	0.61 $\pm$ 0.2
wdbc	0.07 $\pm$ 0.0	0.75 $\pm$ 0.53	2.71 $\pm$ 0.5	0.83 $\pm$ 0.1

and among the results with high proportion of occurrence ( $\geq 0.5$ ). This indicates that, in general, the best results of polynomial-based kernels are achieved by simpler functions, rather than more complex ones. This observation can be counterintuitive and we consider that it justifies an in-depth analysis of the optimal parameters of SVM kernels, which unfortunately is beyond the scope of the present study.

The performance distribution of all Gegenbauer-based classifiers, in terms of accuracy (black crosses) and PSV (red points) is illustrated in Fig. 6 for every dataset; darker shades indicate better classifiers. Notice that the plot for each dataset contains about 52,500 pairs of points, each pair corresponding to an SVM classifier that was trained and evaluated. Altogether the plots contain the performance indexes of more than 1,000,000 classifiers as function of the parameter  $\alpha$ . This plot is a useful tool for the analysis of the SVMs behavior with Gegenbauer kernels.

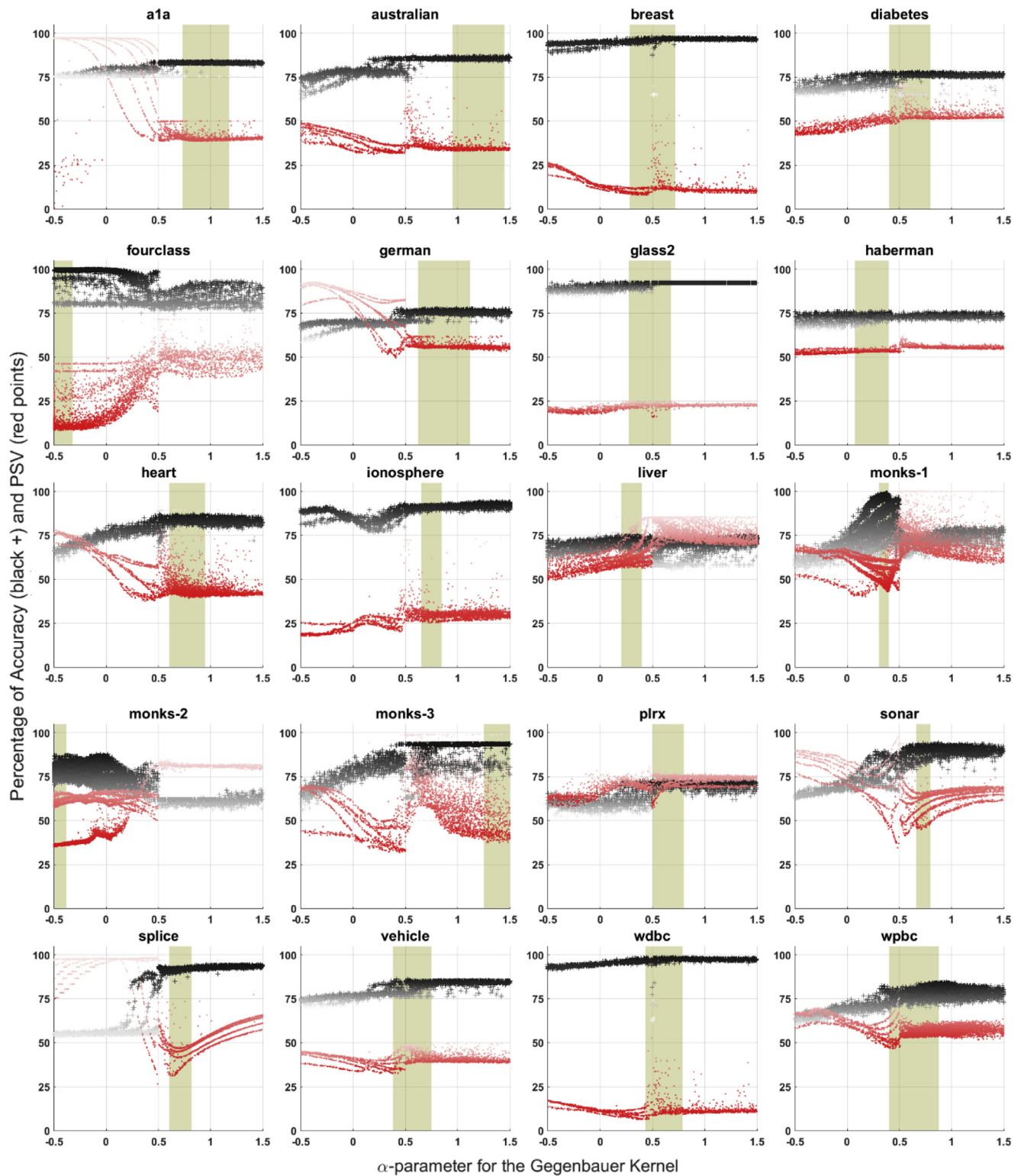
A high accuracy and a small PSV is desirable of any SVM classifier. This combination was found most often for  $\alpha \in [0.3, 0.9]$ . In the plots, the range of  $\alpha$  that produces the best classifiers is indicated by a shaded rectangle. SVM classifiers can suffer from overfitting (PSV 70% or higher); this can be seen in the case of datasets like 'Monks-1', 'Monks-3', and 'Liver' for  $\alpha > 0.5$  and the

**Table 11**

Best values of the polynomial degree (proportion of occurrence within top 25%).

Dataset	MCK	GCK	POLY	LEGEN	CHEBY	s-HERM	GEGEN
a1a	6 (0.20)	1 (0.64)	1 (1.00)	2 (0.85)	4 (0.38)	4 (0.37)	2 (0.19)
Australian	1 (0.96)	1 (0.87)	2 (0.55)	6 (0.31)	1 (0.67)	5 (0.28)	5 (0.24)
Breast	1 (0.31)	1 (0.76)	1 (0.47)	4 (0.44)	3 (0.60)	4 (0.25)	5 (0.21)
Diabetes	1 (0.52)	1 (0.50)	2 (0.48)	2 (0.91)	1 (0.82)	6 (0.28)	2 (0.23)
Fourclass	4 (0.19)	5 (0.25)	5 (1.00)	5 (0.39)	3 (0.26)	6 (0.47)	6 (0.36)
German	1 (0.25)	1 (0.90)	1 (0.38)	3 (0.94)	4 (0.43)	5 (0.30)	1 (0.19)
Glass2	1 (0.28)	4 (0.41)	2 (0.25)	2 (0.44)	6 (0.30)	1 (0.50)	2 (0.24)
Haberman	1 (0.36)	2 (0.45)	4 (0.49)	3 (0.58)	1 (0.85)	6 (0.50)	2 (0.21)
Heart	1 (0.26)	1 (0.43)	1 (0.33)	3 (0.67)	1 (0.69)	6 (0.28)	2 (0.19)
Ionosphere	1 (0.53)	1 (0.70)	2 (0.83)	3 (0.55)	3 (0.48)	5 (0.35)	5 (0.21)
Liver	2 (0.38)	2 (0.39)	2 (0.40)	2 (0.48)	1 (0.75)	2 (0.27)	3 (0.23)
Monks-1	4 (0.25)	1 (0.46)	2 (0.58)	3 (0.53)	3 (0.61)	2 (0.49)	2 (0.43)
Monks-2	2 (0.63)	1 (0.80)	2 (0.98)	2 (1.00)	3 (0.57)	2 (0.54)	1 (0.89)
Monks-3	1 (0.28)	1 (0.62)	1 (0.87)	3 (0.62)	1 (0.46)	6 (0.25)	2 (0.23)
plrx	1 (0.29)	2 (0.43)	2 (0.30)	6 (0.33)	4 (0.27)	1 (0.19)	4 (0.20)
Sonar	1 (0.23)	1 (0.45)	3 (0.33)	2 (1.00)	1 (1.00)	2 (0.32)	5 (0.25)
Splice	1 (0.48)	1 (0.72)	3 (0.46)	2 (0.78)	1 (0.60)	5 (0.66)	2 (0.21)
Vehicle	1 (0.84)	1 (0.67)	3 (0.51)	2 (0.96)	1 (0.50)	6 (0.44)	5 (0.21)
wdbc	1 (0.24)	1 (0.71)	1 (0.43)	2 (0.56)	1 (0.65)	4 (0.25)	2 (0.21)
wdbc	1 (0.38)	1 (0.68)	1 (0.84)	3 (0.43)	3 (0.60)	2 (0.50)	1 (0.37)





**Fig. 6.** Performance of the Gegenbauer spectrum of classifiers obtained by varying its  $\alpha$  parameter. Red dots represent the PSV while black '+' symbols represent the Accuracy. Darker shades indicate classifiers with higher performance metrics. The ranges in which the optimal values of the  $\alpha$ -parameter were found are indicated by a shaded rectangle in the corresponding plot.

dataset 'Splice' for  $\alpha < 0.5$ . For the majority of the datasets, the PSV values exhibit a characteristic behavior: points belonging to the same polynomial degree are grouped together and form distinctive curves which are distributed along the vertical direction. For the accuracy values, this behavior is much less notorious with the exception of a few datasets like 'Fourclass' and 'Splice'.

#### Summary of findings and observations Gegenbauer and s-Hermite:

- Obtained better combinations of accuracy and PSV indexes than the rest of the kernels.

- The hypothesis that these kernels require a smaller PSV than other kernels and maintain high accuracy is confirmed.
- Polynomials of high degrees (e.g. bigger than 6) have higher computational demand compared to lower degrees; our proposed kernels were able to outperform the rest of the polynomial kernels using only small degrees.

Findings related to the parameter  $\alpha$  in the Gegenbauer family:

- $\alpha = 0$ , special case (Chebyshev type-I), obtained the worst results within the family.
- $0.3 \leq \alpha \leq 0.9$ , most of the best Gegenbauer-based classifiers were found in this interval.
- $\alpha > 0.5$ , classifiers in this interval are fairly equivalent to each other in terms of accuracy. PSV becomes a valuable criterion to define the optimum.
- $\alpha < 0.5$ , the polynomials in this interval already have an appropriate amplitude and therefore a scaling and weight functions are unnecessary to avoid the Explosion and Annihilation effects.

## 6. Conclusions and future work

The most important conclusion is that orthogonal polynomial kernels following our proposed formulation are as effective as the benchmark RBF kernel, with the added advantage of being more efficient with respect to the required amount of support vectors. It is our understanding that this behavior is due to two main reasons: first, that the orthogonal polynomials reduce the data representation redundancy in the feature space, leading to less data points being considered as support vectors; second, that our formulation allows the generation of a large amount of new kernels from which an optimum classifier may be obtained.

Another conclusion is that the Gegenbauer kernel family includes Legendre and Chebyshev kernels, with improved performance. Within the Gegenbauer family, many other classifiers were found that overcome these two special cases. Under our formulation, Gegenbauer should be preferred over other orthogonal polynomials because the required scaling function can be analytically computed (contrary, for instance, to Hermite).

Finally, our results suggest that a trade-off between accuracy and PSV indexes must always be considered when working with SVM classifiers, since employing just one of these hides relevant information that can affect the generalization performance of a classifier. Our future work includes exploring the optimization of classifiers based on a combined performance index and further improvements to our orthogonal polynomial kernel formulation.

## Acknowledgment

This work was partially supported by the National Council of Science and Technology of Mexico (CONACYT) through grants: 375524 (Luis C. Padierna) and CATEDRAS-2598 (A. Rojas).

## References

- [1] J. Shawe-Taylor, N. Cristianini, *Kernel Methods for Pattern Analysis*, Cambridge University Press, New York, 2004.
- [2] N. Cristianini, J. Shawe-Taylor, *Kernel Methods for Pattern Analysis*, Cambridge University Press, New York, 2004.
- [3] A. Shigeo, *Support Vector Machines for Pattern Classification*, Springer, New York, 2010.
- [4] N. Deng, Y. Tian, C. Zhang, *Support Vector Machines*, CRC Press, Boca Ratón, 2013.
- [5] J. Shawe-Taylor, S. Sun, A review of optimization methodologies in support vector machines, *Neurocomputing* (2011) 3609–3618.
- [6] V. Vapnik, *Statistical Learning Theory*, New York, John Wiley and Sons, 1998.
- [7] C. Hsu, C. Chang, C. Lin, A practical guide to support vector classification, 2003. [Online]. Available <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>.
- [8] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, C. Watkins, Text classification using string kernels, *J. Mach. Learn. Res.* 2 (2002) 419–444.

- [9] A.D. Essam, T. Hamza, New empirical nonparametric kernels for support vector machines classification, *Appl. Soft Comput.* (13) (2013) 1759–1765.
- [10] N. Ye, R. Sun, Y. Liu, L. Cao, Support vector machine with orthogonal Chebyshev kernel, 18th International Conference on Pattern Recognition, Washington, 2006.
- [11] Z. Pan, H. Chen, X. You, Support vector machine with orthogonal Legendre kernel, *International Conference on Wavelet Analysis and Pattern Recognition*, Xian, 2012.
- [12] S. Ozer, C.H. Chen, Generalized Chebyshev kernels for support vector classification, *ICPR International Conference On Pattern Recognition*, Tampa, FL, USA, 2008.
- [13] S. Ozer, C. Chen, H. Cirpan, A set of new Chebyshev kernel functions for support vector machine pattern classification, *Pattern Recognit.* 44 (7) (2011) 1435–1447.
- [14] H.V. Moghaddam, J. Hamidzadeh, New Hermite orthogonal polynomial kernel and combined kernels in support vector machine classifier, *Pattern Recognit.* 60 (2016) 921–935.
- [15] L. Sun, K.-A. Toh, Z. Lin, A center sliding Bayesian binary classifier adopting orthogonal polynomials, *Pattern Recognit.* 48 (6) (2015) 2013–2028.
- [16] Y. Liu, S. Liao, H. Lin, Y. Yue, W. Wang, Infinite Kernel learning: generalization bounds and algorithms, *Thirty-First AAAI Conference on Artificial Intelligence (AAAI-17)*, 2017.
- [17] O. Chapelle, V. Vapnik, O. Bousquet, S. Mukherjee, Choosing multiple parameters for support vector machines, *Mach. Learn.* 46 (1–3) (2002) 131–159.
- [18] J. Mercer, Functions of positive and negative type, and their connection with the theory of integral equations, in: *Philosophical Transactions of the Royal Society of London, Series A*, 1909, pp. 415–446.
- [19] W. Gautschi, Orthogonal polynomials: applications and computation, *Acta Numer.* (1996) 45–119.
- [20] M. Parodi, J.C. Gómez, Legendre polynomials based feature extraction for on-line signature verification. Consistency analysis of feature combinations, *Pattern Recognit.* 47 (1) (2014) 128–140.
- [21] J. Zhao, G. Yan, B. Feng, W. Mao, J. Bai, An adaptive support vector regression based on a new sequence of unified orthogonal polynomials, *Pattern Recognit.* (2013) 899–913.
- [22] V. Totik, Orthogonal polynomials, *Surv. Approximation Theory* (2005) 70–125.
- [23] C. Dunkl, Y. Xu, *Orthogonal Polynomials of Several Variables*, second ed., Cambridge University Press, Cambridge, 2014.
- [24] T. Chihara, *An Introduction to Orthogonal Polynomials*, Gordon and Breach, New York, 1978.
- [25] Y. Xu, On discrete orthogonal polynomials of several variables, *Adv. Appl. Math.* 33 (2004) 615–632.
- [26] F. Olver, D. Lozier, R. Boisvert, C. Clark, *NIST Handbook of Mathematical Functions*, Cambridge University Press, New York, 2010.
- [27] M. Tian, W. Wang, Some sets of orthogonal polynomial Kernel functions, *Appl. Soft Comput.* 61 (2017) 742–756.
- [28] D.S. Kim, T. Kim, S.-H. Rim, Some identities involving Gegenbauer polynomials, *Adv. Diff. Equat.* 2012 (1) (2012) 219.
- [29] B. Spencer, *The Classical Orthogonal Polynomials* (2015).
- [30] N. Cristianini, J. Shawe-Taylor, *An Introduction to SVM and Other Kernel Based Methods*, Cambridge University Press, Cambridge, U.K, 2000.
- [31] S.I. Valdez, A. Hernández, S. Botello, A Boltzmann based estimation of distribution algorithm, *Inf. Sci.* (2013) 126–137.
- [32] T. Meng, W. Wenjian, Some sets of orthogonal polynomial kernel functions, *Appl. Soft Comput.* 61 (2017) 742–756.
- [33] A. López, X. Li, W. Yu, Support vector machine classification for large datasets using decision tree and fisher linear discriminant, *Future Gener. Comput. Syst.* 36 (2014) 57–65.
- [34] W. Zong, G.-B. Huang, Y. Chen, Weighted extreme learning machine for imbalance learning, *Neurocomputing* 101 (2013) 229–242.
- [35] J. Zhao, Z. Yang, X. Yitian, Nonparallel least square support vector machine for classification, *Appl. Intell.* (2016) 1–10.
- [36] M. Tanveer, Newton method for implicit Lagrangian twin support vector machines, *Int. J. Mach. Learn. Cybern.* 6 (6) (2015) 1029–1040.
- [37] M. Nekkaa, D. Boughaci, A memetic algorithm with support vector machine for feature selection and classification, *Memetic Comput.* 7 (1) (2015) 59–73.
- [38] C.-C. Chang, C.-J. Lin, LIBSVM: a library for support vector machines, *ACM Trans. Intell. Syst. Technol.* 2 (3) (2011) 1–27.
- [39] S. García, A. Fernández, J. Luengo, F. Herrera, Advanced nonparametric tests for multiple comparisons in the design of experiments in computational intelligence and data mining: experimental analysis of power, *Inf. Sci.* 180 (10) (2010) 2044–2064.
- [40] N. Verbiest, K. Vermeulen, A. Teredesai, Evaluation of classification methods, in: *Data Classification: Algorithms and Applications*, CRC Press, Boca Raton, 2015, pp. 633–652.

**Luis Carlos Padierna** received his Ph.D. in Computer Science from the Instituto Tecnológico de León in 2018. His research interests include support vector machines, machine learning, big data and their applications to industry and biomedicine.

**Martín Carpio** received his Ph.D. from the Universidad de Guanajuato, México, in 1995. Currently, he is a Senior Lecturer and Researcher at the Instituto Tecnológico de León. His research interests include computational optimization, simulation processes, statistical data analysis and evolutionary computation.

**Alfonso Rojas-Domínguez** received his Ph.D. from the University of Liverpool, UK, in 2007. Currently, he holds a position as a CONACYT Research Fellow at the Instituto Tecnológico de León. His research interests include automated image processing and analysis, computational intelligence and machine learning.

**Héctor Puga** received his Ph.D. from the Universidad de Guanajuato, México, in 2002. Currently, he leads the Intelligent Systems research group at the Instituto Tecnológico de León. His research interests include optimization, hyper-heuristics and intelligent systems.

**Héctor Fraire** received his Ph.D. in Computer Science from the Centro Nacional de Investigación y Desarrollo Tecnológico, Cuernavaca, México in 2005. His scientific interests include metaheuristic optimization and machine learning.