



**HAL**  
open science

# The Christoffel function for supervised learning: theory and practice

Srećko Đurašinović

► **To cite this version:**

Srećko Đurašinović. The Christoffel function for supervised learning: theory and practice. Machine Learning [stat.ML]. 2022. hal-03768886

**HAL Id: hal-03768886**

**<https://laas.hal.science/hal-03768886>**

Submitted on 5 Sep 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



**Master 2 Applied Mathematics, Statistics**  
Mathematics and Economic Decision

*The Christoffel function for supervised learning: theory  
and practice*  
Internship report

**Student:**

*Srećko Đurašinović*

**Supervisors:**

*Victor Magron*

Full-Time Researcher, LAAS-CNRS

*Jean-Bernard Lasserre*

Emeritus Research Director, LAAS-CNRS

**Keywords:** Christoffel-Darboux Kernel, Polynomial Optimization, Machine Learning, Data Analysis

*July, 2022*



# Summary

Notations . . . . .	4
Introduction . . . . .	5
<b>1 Christoffel function - Selected properties</b>	<b>7</b>
1.1 Preliminaries . . . . .	7
1.2 Reproducing Kernel Hilbert Spaces . . . . .	8
1.3 Christoffel-Darboux kernel . . . . .	9
1.4 Christoffel function . . . . .	11
1.4.1 Some properties . . . . .	13
<b>2 Application in machine learning</b>	<b>17</b>
2.0.1 Introductory remarks . . . . .	17
2.1 Supervised learning . . . . .	17
2.2 Christoffel-Darboux classifier . . . . .	18
2.2.1 Consistency result . . . . .	20
2.3 Empirical evaluations . . . . .	22
2.3.1 <i>Iris</i> data . . . . .	23
2.3.2 Toy data sets . . . . .	24
2.4 Extensions . . . . .	26
Conclusion . . . . .	28
Bibliography . . . . .	29

# Notations

- $\mathcal{M}(K), \mathcal{M}_+(K), \mathcal{P}(K)$  - respectively, set of finite, signed, regular Borel measures, set of positive finite measures, set of probabilities on compact set  $K$
- $\mathcal{C}(K, K'), \mathcal{C}_b(K, K'), \mathcal{C}_0(K, K')$  - respectively, set of continuous, continuous and bounded, continuous vanishing at infinity functions from  $K$  to  $K'$ . If the functions are real-valued, we denote them simply by  $\mathcal{C}(K)$  (resp.  $\mathcal{C}_b(K), \mathcal{C}_0(K)$ )
- $\mathcal{C}(K)^*$  - set of real-valued and continuous on  $\mathcal{C}(K)$
- $\mathbb{R}^{n \times p}$  - set of matrices with  $n \in \mathbb{N} \setminus \{0\}$  lines,  $p \in \mathbb{N} \setminus \{0\}$  columns, with real coefficients/entries
- $\mathbf{I}_n$  - identity matrix of size  $n \times n$ ,  $n \in \mathbb{N} \setminus \{0\}$
- $X^Y$  - space of  $X$ -valued functions defined on  $Y$
- $\text{Lip}_L(X, Y)$  - set of  $L$ -Lipschitz maps from  $X$  to  $Y$
- $A^{\mathbb{N}}$  - set of sequences taking values in the set  $A$
- $B(x, r)$  - closed Euclidean ball of radius  $r > 0$ , centered at  $x$
- $\mathcal{B}(X)$  - Borel sigma-algebra on  $X$
  
- $\langle \cdot, \cdot \rangle_{X \times Y}$  - generic pairing (bi-linear form) between the spaces  $X$  and  $Y$  (the subscript is usually omitted if  $X$  and  $Y$  are easily inferred)
- $\llbracket k, k' \rrbracket$  - interval containing all integers from  $k$  to  $k'$  (included)
- $A \times B$  - Cartesian product of the sets  $A$  and  $B$
- $\|\cdot\|$  - generic norm on some normed vector space
- If  $M$  is a matrix, then  $\|M\|_{\infty} = \sup_{\|x\|_{\infty}=1} \|Mx\|_{\infty}$  is the matrix  $\|\cdot\|_{\infty}$  norm
- $M^T$  - transposition of a matrix  $M$
- $|\cdot|$  - absolute value
- $d_X$  - distance on some metric space  $X$
  
- $\delta_a$  - the Dirac mass concentrated at the point  $a$
- $\delta_{ij}$  - Kronecker delta
- $\mathbf{1}_A$  - Indicator function (equal to 1 if  $A$  is realised, and to 0 otherwise)
- $\wedge, \vee$  - the min and max operators, i.e.  $a \wedge b := \min\{a, b\}$  and  $a \vee b := \max\{a, b\}$
  
- id - the identity map
- 1 - the constant map equal to 1
- 0 - the constant map equal to 0
- $\pi_X$  - projection on the set  $X$
- $\text{diam}(X)$  - diameter of a set  $X$ , i.e.  $\text{diam}(X) := \sup_{x_1, x_2 \in X} d_X(x_1, x_2)$
- $\text{val}(P)$  - optimal value of the optimization program  $(P)$
- $\binom{n}{p} = \frac{n!}{p!(n-p)!}$  - Binomial coefficient
  
- $\mathbb{R}[\mathbf{x}]$  - the ring of real-valued multivariate polynomials in variables  $\mathbf{x} = (x_1, \dots, x_p)$ . The finite dimensional vector space consisting of polynomials of degree at most  $d \in \mathbb{N}$  is denoted by  $\mathbb{R}_d[\mathbf{x}] \subset \mathbb{R}[\mathbf{x}]$ .
- $\mathbb{N}_d^p := \{\boldsymbol{\alpha} \in \mathbb{N}^p \mid |\boldsymbol{\alpha}| := \sum_{i=1}^p \alpha_i \leq d\}$  - the set of multi-indices, with cardinal  $s(d) = \binom{p+d}{p}$ .

# Introduction

Since its inception, machine learning has provided us with many different methods and approaches for handling classification problems. Interested readers can find in [10, 12] a non-exhaustive overview of these methods. In this report, we introduce an additional tool - Christoffel-Darboux kernel - and we exploit its properties, mainly in the context of supervised learning.

Christoffel-Darboux kernel is a very well known tool among the researchers working on the approximation theory. For a long time, its properties were only studied in this context, focusing on the relationship with orthogonal polynomials [9] that can be used for interpolation, .

However, recent works [1, 2, 4, 11, 13] have shed a very different light on this particular kernel by showing its appealing potential in the realm of data science. Indeed, some of the most important features of the Christoffel function can be easily connected to the classical problems in statistics and machine learning: support estimation, outlier detection, graph recovery, free probabilities and many others.

An additional task of equally great importance in data science - supervised learning - can be tackled in this context. We argue that the Christoffel function provides a very intuitive way for solving classification problems, while displaying quite good theoretical properties at the same time [4, 5, 6].

This report will be organised as follows: in the first part, we will introduce the Christoffel function and study its properties that can be useful for solving problems that arise in data science. Then, we will particularly focus on the problem of supervised learning and present some theoretical guarantees that could justify the use of the Christoffel function in this context. Finally, the last part will be devoted to some empirical evaluations and comparisons with some state-of-art methods. Additionally, directions for further research and generalizations will be suggested.



# Chapter 1

## Christoffel function - Selected properties

### 1.1 Preliminaries

Let us start by introducing the notations that will be used throughout this report.

We let  $\mathbf{v}_d(\mathbf{x}) := (\mathbf{x}^\alpha)_{\alpha \in \mathbb{N}_d^p}$  denote the monomial basis of the vector space of  $\mathbb{R}[\mathbf{x}]_d$  (we suppose that the elements of this basis are arranged in the *graded lexicographic order*<sup>1</sup>).

**Definition 1.1.1** (Moment matrix). [7] Let  $\Omega \subset \mathbb{R}^p$  be compact, with non-empty interior. Let  $\phi$  be a Borel measure supported on  $\Omega$ . The moment matrix of order (degree)  $d \in \mathbb{N}$  associated to  $\phi$ , denoted by  $\mathbf{M}_d(\phi)$ , is a  $s(d) \times s(d)$  dimensional real symmetric matrix, whose rows and columns are indexed by monomials in  $\mathbf{v}_d(\mathbf{x})$ , and whose entry  $(\alpha, \beta)$  is given by

$$\mathbf{M}_d(\phi)(\alpha, \beta) := \int_{\Omega} \mathbf{x}^{\alpha+\beta} d\phi, \quad \text{for any } \alpha, \beta \in \mathbb{N}_d^p. \quad (1.1)$$

For the purpose of illustration, with  $p = d = 2$ , we would obtain the following matrix in  $\mathbb{R}^{6 \times 6}$ :

$$\mathbf{M}_2(\phi) = \begin{bmatrix} \int_{\Omega} 1 d\phi & \int_{\Omega} x_1 d\phi & \int_{\Omega} x_2 d\phi & \int_{\Omega} x_1^2 d\phi & \int_{\Omega} x_1 x_2 d\phi & \int_{\Omega} x_2^2 d\phi \\ \cdot & \int_{\Omega} x_1^2 d\phi & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \int_{\Omega} x_2^2 d\phi & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \int_{\Omega} x_1^4 d\phi & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \int_{\Omega} x_1^2 x_2^2 d\phi & \cdot \\ \int_{\Omega} x_2^2 d\phi & \int_{\Omega} x_1 x_2^2 d\phi & \int_{\Omega} x_2^3 d\phi & \int_{\Omega} x_1^2 x_2^2 d\phi & \int_{\Omega} x_1 x_2^3 d\phi & \int_{\Omega} x_2^4 d\phi \end{bmatrix}. \quad (1.2)$$

If we consider  $\mathbf{v}_d(\mathbf{x})$  to be a column vector of  $\mathbb{R}^{s(d)}$ , then we can write the moment matrix in a more compact form as

$$\mathbf{M}_d(\phi) = \int_{\Omega} \mathbf{v}_d(\mathbf{x}) \mathbf{v}_d(\mathbf{x})^T d\phi, \quad (1.3)$$

where the integral should be understood component-wise.

Finally,  $\mathbf{v}_d(\mathbf{x})$  can denote a more general polynomial basis, not necessarily the monomial one. The meaning will always be made clear from the context.

---

<sup>1</sup>To order monomials, we first compares their total degree, i.e. the sum of all their exponents, and in case of a tie, we apply *lexicographic* order. This one, in turn, first compares the exponents of  $x_1$  in the monomials, and if they are the same, then the exponents of  $x_2$  are compared, and so forth.



## 1.2 Reproducing Kernel Hilbert Spaces

In order to better describe the ambient space and the objects that we will manipulate throughout this report, let us briefly recall the idea of Reproducing Kernel Hilbert Spaces (RKHS).

**Definition 1.2.1** (Kernel). Let  $X$  be a set. A kernel is any symmetric map  $k : X \times X \rightarrow \mathbb{R}$ , i.e.  $k(\mathbf{x}, \mathbf{y}) = k(\mathbf{y}, \mathbf{x})$ , for all  $\mathbf{x}, \mathbf{y} \in X$ . If, in addition, the matrix  $(k(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n$  is positive semi-definite for all  $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ , the kernel  $k$  is called positive semi-definite.

Notice that, by definition, a kernel  $k$  is positive semi-definite if  $\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$  for all  $n \in \mathbb{N}$ ,  $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ , and  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ .

**Definition 1.2.2** (RKHS). For  $X \subset \mathbb{R}^p$ , let  $(H, \langle \cdot, \cdot \rangle_H)$  be the Hilbert space of real-valued functions  $f : X \rightarrow \mathbb{R}$  defined on  $X$ . We say that  $H$  is a Reproducing Kernel Hilbert Space if the evaluation functional  $L_{\mathbf{x}}$  is continuous on  $H$ , i.e., if for any  $\mathbf{x} \in X$ , the functional  $H \ni f \mapsto L_{\mathbf{x}}(f) := f(\mathbf{x})$  is continuous.

We can now connect explicitly the notions from the last two definitions.

**Definition 1.2.3** (Reproducing kernel). A kernel  $k$  is called *reproducing kernel* if it satisfies the *reproducing property*:

$$\forall \mathbf{x} \in X, \forall f \in H, L_{\mathbf{x}}(f) = \langle k(\cdot, \mathbf{x}), f \rangle_H = f(\mathbf{x}). \quad (1.4)$$

**Theorem 1.2.1.** [4, 8] For  $X \subset \mathbb{R}^p$ , let  $H$  be the Hilbert space of real-valued functions defined on  $X$ . Let  $k : X \times X \rightarrow \mathbb{R}$  be a kernel.

1. If  $H$  is a RKHS, then it admits a positive semi-definite reproducing kernel.
2. If  $k$  is positive semi-definite, there exists (a unique) RKHS with  $k$  as its reproducing kernel.

*Proof.* Let  $X, H$  and  $k$  be defined as in the statement of the theorem.

1. Let  $\mathbf{x}, \mathbf{y} \in X$ . Since  $H$  is RKHS, then, by definition, the linear functional  $L_{\mathbf{x}}$  is continuous. Hence, the classical Riesz representation theorem ensures the there exists some  $k_{\mathbf{x}} \in H$  such that  $L_{\mathbf{x}}(f) = \langle f, k_{\mathbf{x}} \rangle_H$  for all  $f \in H$ . Using the same arguments, there exists  $k_{\mathbf{y}} \in H$  such that  $L_{\mathbf{y}}(f) = \langle f, k_{\mathbf{y}} \rangle_H$  for all  $f \in H$ . Let us define on the product space  $X \times X$  the mapping  $k$  by  $k(\mathbf{x}, \mathbf{y}) = \langle k_{\mathbf{x}}, k_{\mathbf{y}} \rangle_H$ . By construction,  $k$  satisfies the reproducing property. Moreover, for any  $n \in \mathbb{N}$ ,  $\mathbf{x}_1, \dots, \mathbf{x}_n \in X$ , and  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ , we have

$$\sum_{i,j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{i,j=1}^n \alpha_i \alpha_j \langle k_{\mathbf{x}_i}, k_{\mathbf{x}_j} \rangle_H \quad (1.5)$$

$$= \left\langle \sum_{i=1}^n \alpha_i k_{\mathbf{x}_i}, \sum_{i=1}^n \alpha_i k_{\mathbf{x}_i} \right\rangle_H = \left\| \sum_{i=1}^n \alpha_i k_{\mathbf{x}_i} \right\|_H^2 \geq 0, \quad (1.6)$$

proving that  $k$  is positive semi-definite.

2. Consider the mapping  $\phi : X \rightarrow \mathbb{R}^X$  such that  $\phi(\mathbf{x}) = k(\cdot, \mathbf{x})$  for all  $\mathbf{x} \in X$ . Then, we let  $V := \text{vect}(\{\phi(\mathbf{x}) \mid \mathbf{x} \in X\}) = \{f(\cdot) = \sum_{i=1}^n \alpha_i k(\cdot, \mathbf{x}_i) \mid n \in \mathbb{N}, \alpha_i \in \mathbb{R}, \mathbf{x}_i \in X\}$ . Moreover,

we can define for two elements in  $V$ , say  $f(\cdot) = \sum_{i=1}^{n_f} \alpha_i k(\cdot, \mathbf{x}_i)$  and  $g(\cdot) = \sum_{i=1}^{n_g} \beta_i k(\cdot, \mathbf{y}_i)$  the following form:

$$\langle f, g \rangle = \sum_{i=1}^{n_f} \sum_{j=1}^{n_g} \alpha_i \beta_j k(\mathbf{x}_i, \mathbf{y}_j).$$

The assumptions about  $k$  ensure that  $\langle \cdot, \cdot \rangle$  is a valid scalar product. The closure of  $V$  would then form a RKHS, since for any  $f \in \text{closure}(V)$ ,

$$|f(\mathbf{x})| = |L_{\mathbf{x}}(f)| = |\langle f, k(\cdot, \mathbf{x}) \rangle| \leq \|k(\cdot, \mathbf{x})\|_H \|f\|_H, \quad (1.7)$$

where the last inequality follows from the Cauchy-Schwartz inequality. □

Thus, we have managed to obtain a Hilbert space satisfying some additional reproducing properties. This additional space structure can also be obtained with Christoffel-Darboux kernels that we will introduce in a sequel. It turns out that these features can be exploited in the context of machine learning.

### 1.3 Christoffel-Darboux kernel

In the spirit of the previous section, one may see the space  $\mathbb{R}_d[\mathbf{x}]$  (which is a subspace of  $L^2(\phi)$ ) as the RKHS.

If, for instance, the measure  $\phi$  is compactly supported and absolutely continuous with respect to the Lebesgue measure, then we can show that the following bi-linear form

$$\langle \cdot, \cdot \rangle_{\phi}: \mathbb{R}_d[\mathbf{x}] \times \mathbb{R}_d[\mathbf{x}] \longrightarrow \mathbb{R} \quad (1.8)$$

$$(P, Q) \longmapsto \langle P, Q \rangle_{\phi} := \int_{\Omega} PQ d\phi \quad (1.9)$$

defines a valid scalar product, making  $(\mathbb{R}_d[\mathbf{x}], \langle \cdot, \cdot \rangle_{\phi})$  a finite-dimensional Hilbert space of polynomial functions from  $\mathbb{R}^p$  to  $\mathbb{R}$  [4]. Since every linear map defined on a finite-dimensional space is continuous, we deduce that  $(\mathbb{R}_d[\mathbf{x}], \langle \cdot, \cdot \rangle_{\phi})$  is indeed a RKHS. Its associated reproducing kernel is known as the Christoffel-Darboux kernel. In order to define it properly, let us recall the next definition:

**Definition 1.3.1.** Let  $(P_{\alpha})_{\alpha \in \mathbb{N}^p}$  be a family of polynomials in  $\mathbb{R}[\mathbf{x}]$ . We say that this family is *orthonormal with respect to  $\phi$*  if

$$\int_{\Omega} P_{\alpha}(\mathbf{x}) P_{\beta}(\mathbf{x}) d\phi(\mathbf{x}) = \mathbf{1}_{\alpha=\beta}, \quad \forall \alpha, \beta \in \mathbb{N}^p. \quad (1.10)$$

**Definition 1.3.2** (Christoffel-Darboux kernel). The *Christoffel-Darboux kernel* is the reproducing kernel associated to  $(\mathbb{R}_d[\mathbf{x}], \langle \cdot, \cdot \rangle_{\phi})$ . It can be defined via the following relationship

$$K_d^{\phi}(\mathbf{x}, \mathbf{y}) := \sum_{i=1}^{s(d)} P_i(\mathbf{x}) P_i(\mathbf{y}), \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^p, \quad (1.11)$$

where  $(P_i)_{i=1}^{s(d)}$  is an orthonormal basis of  $\mathbb{R}_d[\mathbf{x}]$ .

*Remark 1.3.1.* The kernel  $K_d^\phi$  satisfies the reproducing property. Indeed, if  $P \in \mathbb{R}_d[\mathbf{x}]$ , then one can write  $P = \sum_{i=1}^{s(d)} p_i P_i$ , for some real numbers  $(p_i)_{i=1}^{s(d)}$ , implying that

$$\forall \mathbf{x} \in X, L_{\mathbf{x}}(P) = \langle K_d^\phi(\cdot, \mathbf{x}), P \rangle_\phi \quad (1.12)$$

$$= \int_{\Omega} P(\mathbf{y}) K_d^\phi(\mathbf{y}, \mathbf{x}) d\phi(\mathbf{y}) = \int_{\Omega} P(\mathbf{y}) \sum_{i=1}^{s(d)} P_i(\mathbf{x}) P_i(\mathbf{y}) d\phi(\mathbf{y}) \quad (1.13)$$

$$= \sum_{i=1}^{s(d)} P_i(\mathbf{x}) \int_{\Omega} P(\mathbf{y}) P_i(\mathbf{y}) d\phi(\mathbf{y}) = \sum_{i=1}^{s(d)} P_i(\mathbf{x}) \sum_{j=1}^{s(d)} p_j \underbrace{\int_{\Omega} P_j(\mathbf{y}) P_i(\mathbf{y}) d\phi(\mathbf{y})}_{\mathbf{1}_{j=i}} \quad (1.14)$$

$$= \sum_{i=1}^{s(d)} p_i P_i(\mathbf{x}) = P(\mathbf{x}). \quad (1.15)$$

An interesting feature of this kernel is that it can be explicitly computed from the moments, i.e., from the moment matrix. Indeed, if we let  $\mathbf{v}_d(\mathbf{x}) := (P_1(\mathbf{x}), \dots, P_{s(d)}(\mathbf{x})) \in \mathbb{R}^{s(d)}$  for some basis  $(P_i)_{i=1}^{s(d)}$  of  $\mathbb{R}_d[\mathbf{x}]$  (not necessarily orthonormal), then for any polynomial  $P \in \mathbb{R}_d[\mathbf{x}]$  there exists some  $\mathbf{p} \in \mathbb{R}^{s(d)}$  such that  $P(\mathbf{x}) = \mathbf{p}^T \mathbf{v}_d(\mathbf{x})$  for any  $\mathbf{x} \in \mathbb{R}^p$ .

In particular, since  $\langle \cdot, \cdot \rangle_\phi$  is a valid scalar product, then for any  $P \neq 0$ , we have

$$\langle P, P \rangle_\phi = \int_{\Omega} \mathbf{p}^T \mathbf{v}_d(\mathbf{x}) (\mathbf{p}^T \mathbf{v}_d(\mathbf{x}))^T d\phi(\mathbf{x}) = \mathbf{p}^T \left( \int_{\Omega} \mathbf{v}_d(\mathbf{x}) \mathbf{v}_d(\mathbf{x})^T d\phi(\mathbf{x}) \right) \mathbf{p} = \mathbf{p}^T \mathbf{M}_d(\phi) \mathbf{p} > 0, \quad (1.16)$$

implying that  $\mathbf{M}_d(\phi) \succ 0$ . This allows us to derive another expression for the Christoffel-Darboux kernel, based on the moments. Namely,

$$K_d^\phi(\mathbf{x}, \mathbf{y}) := \mathbf{v}_d(\mathbf{x})^T \mathbf{M}_d(\phi)^{-1} \mathbf{v}_d(\mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p. \quad (1.17)$$

That this new formulation satisfies the reproducing property is guaranteed by the uniqueness of the reproducing kernel, but it can also be seen from the following computation:  $\forall \mathbf{x} \in X, \forall \mathbf{p} \in \mathbb{R}^{s(d)}$ ,

$$\langle K_d^\phi(\cdot, \mathbf{x}), P \rangle_\phi = \int_{\Omega} P(\mathbf{y}) K_d^\phi(\mathbf{y}, \mathbf{x}) d\phi(\mathbf{y}) = \int_{\Omega} \mathbf{p}^T \mathbf{v}_d(\mathbf{y}) \mathbf{v}_d(\mathbf{y})^T \mathbf{M}_d(\phi)^{-1} \mathbf{v}_d(\mathbf{x}) d\phi(\mathbf{y}) \quad (1.18)$$

$$= \mathbf{p}^T \underbrace{\int_{\Omega} \mathbf{v}_d(\mathbf{y}) \mathbf{v}_d(\mathbf{y})^T d\phi(\mathbf{y})}_{\mathbf{I}_{s(d)}} \mathbf{M}_d(\phi)^{-1} \mathbf{v}_d(\mathbf{x}) \quad (1.19)$$

$$= \mathbf{p}^T \mathbf{v}_d(\mathbf{x}) = P(\mathbf{x}). \quad (1.20)$$

Finally, we can observe that  $\mathbf{M}_d(\phi) = (\langle P_i, P_j \rangle_\phi)_{i,j=1}^{s(d)}$ , so that, if the polynomial basis is orthonormal, we obtain  $\mathbf{M}_d(\phi) = \mathbf{I}_{s(d)}$ , which shows explicitly that expressions (1.11) and (1.17) coincide.

Finally, let us mention that changing the polynomial basis does not change the Christoffel-Darboux kernel, which is summarized in the next proposition.

**Proposition 1.3.1.** Let  $\mathbf{v}_d(\cdot)$  and  $\mathbf{w}_d(\cdot)$  be two different polynomial bases of  $\mathbb{R}_d[\mathbf{x}]$ . Let  $\mathbf{M}_d(\phi)$  and  $\mathbf{N}_d(\phi)$  be defined as in (1.3), with  $\mathbf{v}_d(\cdot)$  and  $\mathbf{w}_d(\cdot)$  respectively. Then, the associated Christoffel-Darboux kernels are equal.

*Proof.* Since  $\mathbf{v}_d(\cdot)$  and  $\mathbf{w}_d(\cdot)$  are two different basis of the vector space  $\mathbb{R}_d[\mathbf{x}]$ , there exists a positive definite matrix  $\Pi \in \mathbb{R}^{s(d) \times s(d)}$  such that  $\mathbf{w}_d(\cdot) = \Pi \mathbf{v}_d(\cdot)$ . Let  $\mathbf{x} \in \mathbb{R}^p$ . Then,

$$\mathbf{w}_d(\mathbf{x})^T (\mathbf{N}_d(\phi))^{-1} \mathbf{w}_d(\mathbf{x}) = \mathbf{w}_d(\mathbf{x})^T \left( \int_{\Omega} \mathbf{w}_d(\mathbf{z}) \mathbf{w}_d(\mathbf{z})^T d\phi(\mathbf{z}) \right)^{-1} \mathbf{w}_d(\mathbf{x}) \quad (1.21)$$

$$= \mathbf{v}_d(\mathbf{x})^T \Pi^T \left( \int_{\Omega} \Pi \mathbf{v}_d(\mathbf{z}) \mathbf{v}_d(\mathbf{z})^T \Pi^T d\phi(\mathbf{z}) \right)^{-1} \Pi \mathbf{v}_d(\mathbf{x}) \quad (1.22)$$

$$= \mathbf{v}_d(\mathbf{x})^T \Pi^T \left( \Pi \int_{\Omega} \mathbf{v}_d(\mathbf{z}) \mathbf{v}_d(\mathbf{z})^T d\phi(\mathbf{z}) \Pi^T \right)^{-1} \Pi \mathbf{v}_d(\mathbf{x}) \quad (1.23)$$

$$= \mathbf{v}_d(\mathbf{x})^T \Pi^T \left( \Pi \mathbf{M}_d(\phi) \Pi^T \right)^{-1} \Pi \mathbf{v}_d(\mathbf{x}) \quad (1.24)$$

$$= \mathbf{v}_d(\mathbf{x})^T \Pi^T (\Pi^T)^{-1} (\mathbf{M}_d(\phi))^{-1} \Pi^{-1} \Pi \mathbf{v}_d(\mathbf{x}) \quad (1.25)$$

$$= \mathbf{v}_d(\mathbf{x})^T (\mathbf{M}_d(\phi))^{-1} \mathbf{v}_d(\mathbf{x}). \quad (1.26)$$

□

What seems to distinguish the Christoffel-Darboux kernel from other kernels commonly used in practice is its profound connection with an underlying measure  $\phi$ . Indeed, this measure is inducing an inner product on the ambient Reproducing Kernel Hilbert Space, which is suitable for many problems in statistics. Other commonly used kernels (polynomial kernel) are not induced by a positive measure.

## 1.4 Christoffel function

In the previous section, the Christoffel-Darboux kernel was introduced. Using the same notations, we can define the Christoffel function, the very central object of this report.

**Definition 1.4.1** (Christoffel function). The *Christoffel function* of order (degree)  $d \in \mathbb{N}$ , denoted by  $\Lambda_d^\phi$ , is defined via the following relationship

$$\Lambda_d^\phi: \mathbb{R}^p \longrightarrow \mathbb{R}_+ \quad (1.27)$$

$$\mathbf{x} \longmapsto \Lambda_d^\phi(\mathbf{x}) := \min_{P \in \mathbb{R}_d[\mathbf{x}]} \left\{ \int_{\Omega} P^2(\mathbf{z}) d\phi(\mathbf{z}), P(\mathbf{x}) = 1 \right\}. \quad (1.28)$$

In other words, to any point  $\mathbf{x} \in \mathbb{R}^p$ , the Christoffel function associates the minimal squared  $\|\cdot\|_\phi$ -norm among the polynomials whose value at  $\mathbf{x}$  is equal to one. This constrained optimization problem is in fact equivalent to a particular quadratic programming problem. Indeed, recalling that any polynomial  $P \in \mathbb{R}_d[\mathbf{x}]$  can be written as  $\mathbf{p}^T \mathbf{v}_d(\cdot)$ , we deduce that the objective function becomes  $\int_{\Omega} \mathbf{p}^T \mathbf{v}_d(\mathbf{z}) \mathbf{v}_d(\mathbf{z})^T \mathbf{p} d\phi(\mathbf{z}) = \mathbf{p}^T \mathbf{M}_d(\phi) \mathbf{p}$ , so that

$$\Lambda_d^\phi(\mathbf{x}) := \min_{\mathbf{p} \in \mathbb{R}^{s(d)}} \left\{ \mathbf{p}^T \mathbf{M}_d(\phi) \mathbf{p}, \mathbf{p}^T \mathbf{v}_d(\mathbf{x}) = 1 \right\}. \quad (1.29)$$

Hence, evaluating the Christoffel function at a particular point is the same as solving a specific convex quadratic programming. The interest of this formulation lies in the fact that this family of optimization problems can be handled very efficiently, even in large dimensions, by some numerical solvers, which is crucial in practical computations. However, the drawback is that for each point, one has to solve a completely new quadratic programming, i.e., no closed formula is available.

Let us consider an alternative way of defining this function:

**Proposition 1.4.1.** The Christoffel function of order (degree)  $d \in \mathbb{N}$ , defined in (1.28), satisfies the following relationship

$$\Lambda_d^\phi(\mathbf{x}) = \frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})} = \frac{1}{\mathbf{v}_d(\mathbf{x})^T \mathbf{M}_d(\phi)^{-1} \mathbf{v}_d(\mathbf{x})}. \quad (1.30)$$

*Proof.* Let  $\mathbf{x} \in \mathbb{R}^p$ . First, we can see that the quotient in (1.30) is well defined thanks to the orthogonal polynomials definition of the Christoffel-Darboux kernel (1.11). Indeed, we have  $K_d^\phi(\mathbf{x}, \mathbf{x}) = \sum_{i=1}^{s(d)} P_i(\mathbf{x})^2 > 0$ , for all  $\mathbf{x} \in \mathbb{R}^p$ .

Let us now consider the polynomial  $P : \mathbf{z} \mapsto \frac{K_d^\phi(\mathbf{z}, \mathbf{x})}{K_d^\phi(\mathbf{x}, \mathbf{x})}$ . By definition,  $P \in \mathbb{R}_d[\mathbf{x}]$ , and since  $P(\mathbf{x}) = \frac{K_d^\phi(\mathbf{x}, \mathbf{x})}{K_d^\phi(\mathbf{x}, \mathbf{x})} = 1$ , we deduce that  $P$  is admissible for the optimization problem in (1.28). This means that

$$\Lambda_d^\phi(\mathbf{x}) \leq \int_{\Omega} \frac{K_d^\phi(\mathbf{z}, \mathbf{x})^2}{K_d^\phi(\mathbf{x}, \mathbf{x})^2} d\phi(\mathbf{z}) = \frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})^2} \int_{\Omega} K_d^\phi(\mathbf{z}, \mathbf{x}) K_d^\phi(\mathbf{z}, \mathbf{x}) d\phi(\mathbf{z}) \quad (1.31)$$

$$= \frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})^2} K_d^\phi(\mathbf{x}, \mathbf{x}) \quad (\text{reproducing property}) \quad (1.32)$$

$$= \frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})}. \quad (1.33)$$

On the other hand, if  $P$  is an admissible solution, then

$$1 = P(\mathbf{x})^2 = \left( \int_{\Omega} K_d^\phi(\mathbf{z}, \mathbf{x}) P(\mathbf{z}) d\phi(\mathbf{z}) \right)^2 \quad (\text{reproducing property}) \quad (1.34)$$

$$\leq \int_{\Omega} K_d^\phi(\mathbf{z}, \mathbf{x})^2 d\phi(\mathbf{z}) \int_{\Omega} P(\mathbf{z})^2 d\phi(\mathbf{z}) \quad (\text{Cauchy-Schwartz}) \quad (1.35)$$

$$= K_d^\phi(\mathbf{x}, \mathbf{x}) \int_{\Omega} P(\mathbf{z})^2 d\phi(\mathbf{z}) \quad (\text{reproducing property}) \quad (1.36)$$

implying that  $\frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})} \leq \Lambda_d^\phi(\mathbf{x})$ . Combining the two inequalities, we deduce that  $\Lambda_d^\phi(\mathbf{x}) = \frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})}$ . Moreover, from computations leading to (1.33), we deduce that the polynomial  $P$  reaches the optimal value  $\frac{1}{K_d^\phi(\mathbf{x}, \mathbf{x})}$ , so that it actually solves the problem in (1.28).  $\square$

Thus, evaluating the Christoffel function at some point  $\mathbf{x} \in \mathbb{R}^p$  can be done by taking the reciprocal of the value of the Christoffel-Darboux kernel along the diagonal  $\mathbf{y} = \mathbf{x}$ . Indeed, recall that  $\mathbf{v}_d(\cdot)$  denotes a monomial basis of  $\mathbb{R}_d[\mathbf{x}]$ , so that  $K_d^\phi(\mathbf{x}, \mathbf{x}) = \mathbf{v}_d(\mathbf{x})^T \mathbf{M}_d(\phi)^{-1} \mathbf{v}_d(\mathbf{x})$  is nothing but a  $2d$ -degree polynomial evaluated at  $\mathbf{x} \in \mathbb{R}^p$ .

An advantage of this formulation is that it gives us an explicit closed-form formula for the Christoffel function, directly related to the Christoffel-Darboux kernel. However, this requires inverting an  $s(d) \times s(d)$ -dimensional matrix, which can be extremely costly in the high dimensional context, i.e., if  $p$  is high, or very unstable numerically, if  $d$  is high, for example.

This function has been known for a very long time among the researchers working on the approximation theory and orthogonal polynomials. The next figure motivates the use of the Christoffel function in statistical applications:

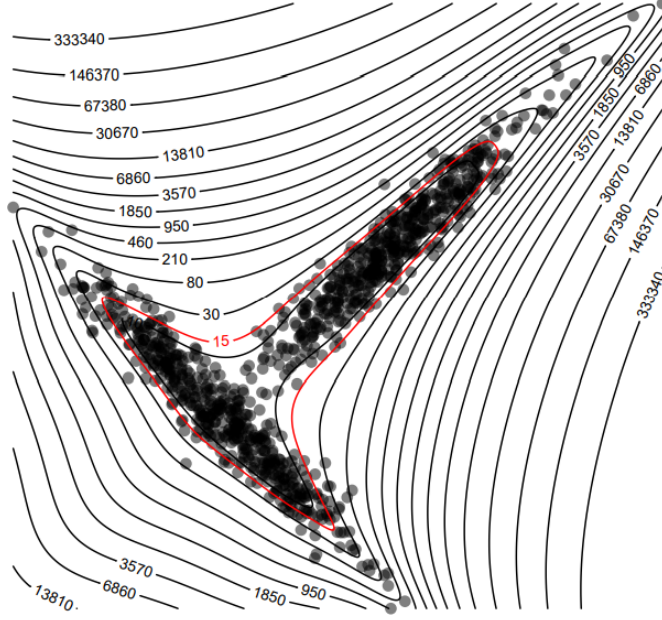


Figure 1.1: Level sets of a Christoffel function associated to the empirical measure [4, page 2].

The Figure 1.1. is obtained from an empirical measure  $\phi_N := \sum_{i=1}^N \delta_{\mathbf{x}_i}$  associated to the cloud of data points  $(\mathbf{x}_i)_{i=1}^N$  in  $\mathbb{R}^2$ , with  $N = 1000$ . The red curve is just one particular level set associated to  $\Lambda_4^{\phi_N}$ . An immediate observation is that these level sets seem to capture very well the shape of the data points, i.e. the support of the underlying measure that the data was generated from. Although this property is true for more general measures and can be obtained from the true theoretical moments, the fact that it can be applied to the empirical measure, which is one of the central objects in statistics, increases the attractiveness of this function.

### 1.4.1 Some properties

Let us start by showing that the Christoffel function can identify pure points asymptotically.

**Proposition 1.4.2.** [3, 4] Let  $\phi \in \mathcal{P}(\Omega)$ , with  $\Omega \subset \mathbb{R}^p$  compact. Let  $\mathbf{x} \in \mathbb{R}^p$ . Then,

$$\lim_{d \rightarrow +\infty} \Lambda_d^\phi(\mathbf{x}) = \phi(\{\mathbf{x}\}). \quad (1.37)$$

*Proof.* Let  $\mathbf{x} \in \mathbb{R}^p$ . It is clear from the variational formulation in (1.28) that the limit in (1.37) exists, since the sequence  $(\Lambda_d^\phi(\mathbf{x}))_{d \in \mathbb{N}}$  is bounded from below by zero and non-increasing. Furthermore, using the properties of an integral with respect to a measure, we deduce that for any admissible  $P \in \mathbb{R}_d[\mathbf{x}]$ ,

$$\int_{\Omega} P(\mathbf{z})^2 d\phi(\mathbf{z}) \geq P(\mathbf{x})\phi(\{\mathbf{x}\}) = \phi(\{\mathbf{x}\}). \quad (1.38)$$

By taking the infimum and then the limit when  $d \rightarrow +\infty$ , we get that  $\lim_{d \rightarrow +\infty} \Lambda_d^\phi(\mathbf{x}) \geq \phi(\{\mathbf{x}\})$ . On the other hand, for any  $d \in \mathbb{N}$ , the polynomial  $P : \mathbf{z} \mapsto (1 - \|\mathbf{z} - \mathbf{x}\|^2)^d$  belongs to  $\mathbb{R}_{2d}[\mathbf{x}]$  and

satisfies  $P(\mathbf{x}) = 1$  so that

$$\Lambda_{2d+1}^\phi(\mathbf{x}) \leq \Lambda_{2d}^\phi(\mathbf{x}) \leq \int_{\Omega} (1 - \|\mathbf{z} - \mathbf{x}\|^2)^{2d} d\phi(\mathbf{z}) \quad (1.39)$$

$$\leq \int_{B(\mathbf{x}, d^{-\frac{1}{4}})} d\phi(\mathbf{z}) + \int_{\Omega \setminus B(\mathbf{x}, d^{-\frac{1}{4}})} (1 - \|\mathbf{z} - \mathbf{x}\|^2)^{2d} d\phi(\mathbf{z}) \quad (1.40)$$

$$\leq \int_{B(\mathbf{x}, d^{-\frac{1}{4}})} d\phi(\mathbf{z}) + \int_{\Omega \setminus B(\mathbf{x}, d^{-\frac{1}{4}})} (1 - (d^{-\frac{1}{4}})^2)^{2d} d\phi(\mathbf{z}) \quad (1.41)$$

$$= \underbrace{\int_{B(\mathbf{x}, d^{-\frac{1}{4}})} d\phi(\mathbf{z})}_{\xrightarrow{d \rightarrow +\infty} \phi(\{\mathbf{x}\})} + \underbrace{\left(1 - \frac{1}{\sqrt{d}}\right)^{2d}}_{\xrightarrow{d \rightarrow +\infty} 0}. \quad (1.42)$$

Therefore, we conclude that  $\lim_{d \rightarrow +\infty} \Lambda_d^\phi(\mathbf{x}) = \phi(\{\mathbf{x}\})$ .  $\square$

The next result we present is related to the average value taken by the Christoffel function inside the support of its associated measure.

**Proposition 1.4.3.** [3, 4] Let  $\Omega \subset \mathbb{R}^p$  be compact and  $\phi \in \mathcal{P}(\Omega)$  be an absolutely continuous probability measure. Let  $X$  be a random variable distributed according to  $\phi$ . Let  $d \in \mathbb{N}$ . Then,

$$\mathbb{E}_\phi \left[ \frac{1}{\Lambda_d^\phi(X)} \right] = s(d). \quad (1.43)$$

*Proof.* Let  $P_1, \dots, P_{s(d)}$  be an orthonormal basis of  $\mathbb{R}[\mathbf{x}]$ . By definition of an expectation,

$$\mathbb{E}_\phi \left[ \frac{1}{\Lambda_d^\phi(X)} \right] = \mathbb{E}_\phi [K_d^\phi(X, X)] = \int_{\Omega} K_d^\phi(\mathbf{x}, \mathbf{x}) d\phi(\mathbf{x}) \quad (1.44)$$

$$= \int_{\Omega} \sum_{i=1}^{s(d)} P_i(\mathbf{x})^2 d\phi(\mathbf{x}) = \sum_{i=1}^{s(d)} \underbrace{\int_{\Omega} P_i(\mathbf{x})^2 d\phi(\mathbf{x})}_{= \|P_i\|_\phi^2 = 1} \quad (1.45)$$

$$= s(d). \quad (1.46)$$

$\square$

Thus, for an absolutely continuous measure on a compact support, the value of the Christoffel function decreases polynomially fast towards zero when we let  $d \rightarrow +\infty$  (recall that  $s(d) = \binom{p+d}{p} \sim d^p$ ).

**Proposition 1.4.4.** [4] Let  $d \in \mathbb{N}$ , and let  $\phi \in \mathcal{P}(\Omega)$ , with  $\Omega \subset \mathbb{R}^p$  compact. Then, for any  $\mathbf{x}_0 \in \mathbb{R}^p$  verifying  $\text{dist}(\mathbf{x}_0, \Omega) \geq \delta > 0$ , one has

$$K_d^\phi(\mathbf{x}_0, \mathbf{x}_0) \geq 2^{\frac{d\delta}{\delta + \text{diam}(\Omega)} - 3}. \quad (1.47)$$

In order to prove the previous proposition, we use the following lemma:

**Lemma 1.4.1.** [9] Let  $d \in \mathbb{N}^*$  and  $\delta \in ]0, 1[$ . Then, there exists  $Q \in \mathbb{R}_{2d}[\mathbf{x}]$  satisfying:

$$Q(\mathbf{0}) = 1, \quad \|\mathbf{x}\| \leq 1 \implies |Q(\mathbf{x})| \leq 1, \quad \text{and } 0 < \delta \leq \|\mathbf{x}\| \leq 1 \implies |Q(\mathbf{x})| \leq 2^{1-\delta d}. \quad (1.48)$$

*Proof of the Proposition 1.4.4.* Let  $d \in \mathbb{N}^*$  and let  $\delta \in ]0, 1[$ . Set  $\tilde{\delta} = \frac{\delta}{\delta + \text{diam}(\Omega)} \in ]0, 1[$ . Consider  $P \in \mathbb{R}_{2d}[\mathbf{x}]$  such that  $P(\mathbf{x}) = Q\left(\frac{\mathbf{x} - \mathbf{x}_0}{\delta + \text{diam}(\Omega)}\right)$ , where  $Q \in \mathbb{R}_{2d}[\mathbf{x}]$  is the polynomial whose existence is assured by Lemma 1.4.1. Since  $P(\mathbf{x}_0) = Q(\mathbf{0}) = 1$ , we deduce that  $P$  is an admissible candidate in the optimization problem (1.28) defining  $\Lambda_{2d}^\phi(\mathbf{x}_0)$ . Moreover, since  $\text{dist}(\mathbf{x}_0, \Omega) \geq \delta$ , we deduce  $\tilde{\delta} \leq \left\| \frac{\mathbf{x} - \mathbf{x}_0}{\delta + \text{diam}(\Omega)} \right\| \leq 1$  for all  $\mathbf{x} \in \Omega$ . Hence, we obtain the following:

$$\Lambda_{2d}^\phi(\mathbf{x}_0) \leq \int P^2(\mathbf{x}) d\phi(\mathbf{x}) \leq \int (2^{1-\tilde{\delta}d})^2 d\phi(\mathbf{x}) = 2^{2-2\tilde{\delta}d} \leq 2^{3-\tilde{\delta}2d}. \quad (1.49)$$

Additionally, since  $\mathbb{R}_{2d}[\mathbf{x}] \subset \mathbb{R}_{2d+1}[\mathbf{x}]$ , we deduce that  $\Lambda_{2d+1}^\phi(\mathbf{x}_0) \leq \Lambda_{2d}^\phi(\mathbf{x}_0)$ . Combining the fact that  $\tilde{\delta} \leq 1$  with the equation (1.49), we deduce that

$$\Lambda_{2d+1}^\phi(\mathbf{x}_0) \leq 2^{3-\tilde{\delta}(2d+1)}. \quad (1.50)$$

Since we have shown the validity of the formula for odd and even degrees separately, we deduce that for any  $d \in \mathbb{N}^*$ ,  $\Lambda_d^\phi(\mathbf{x}_0) \leq 2^{3-\tilde{\delta}d}$ . Recalling the definition of  $\tilde{\delta}$  and using the Proposition 1.4.1, we obtain

$$\Lambda_d^\phi(\mathbf{x}_0) \leq 2^{3-\tilde{\delta}d} \implies K_d^\phi(\mathbf{x}_0, \mathbf{x}_0) \geq 2^{\tilde{\delta}d-3} = 2^{\frac{d\delta}{\delta + \text{diam}(\Omega)}-3}. \quad (1.51)$$

□

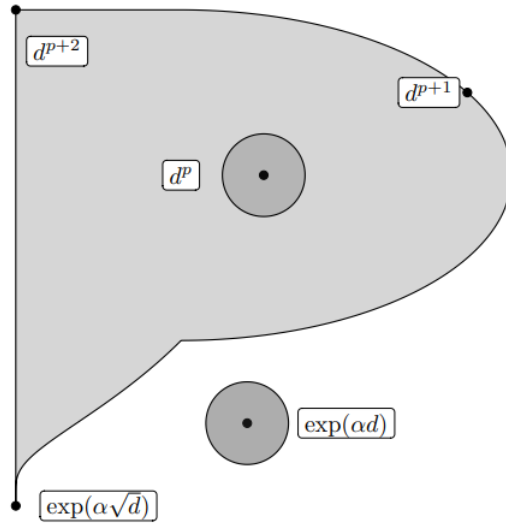


Figure 1.2: (Adapted (link) with permission) Different growth rates for  $(\Lambda_d^\phi(\mathbf{x}))^{-1} = K_d^\phi(\mathbf{x}, \mathbf{x})$ .

The Figure 1.2, together with the Proposition 1.4.3. and the Proposition 1.4.4. provides a rationale for the use of the Christoffel function in data analysis. Indeed, there exists some kind of dichotomy that is manifested through the asymptotically exponentially fast decay of the Christoffel function for the points not belonging to the support of its associated measure. On the other hand, as  $d$  increases, this decay is at most polynomial for the points inside the support. This is essentially the main mechanism through which the support of the underlying measure is identified. Moreover, we can also observe that different rates can be obtained for the points lying on the boundary, depending on its smoothness.



Finally, let us conclude this section by stating a theorem that relates the Christoffel functions associated to a certain measure and its empirical counterpart.

**Theorem 1.4.1.** [2, 4] Let  $d \in \mathbb{N} \setminus \{0\}$ , and let  $\phi \in \mathcal{P}(\Omega)$ , with  $\Omega \subset \mathbb{R}^p$  compact. Let  $n \in \mathbb{N}$  and  $(\mathbf{x}_i)_{i=1}^n$  be a sample of points in  $\mathbb{R}^p$  drawn according to  $\phi$ , and  $\phi_n := \sum_{i=1}^n \delta_{\mathbf{x}_i}$  its associated empirical measure. Then,

$$\|\Lambda_d^{\phi_n} - \Lambda_d^\phi\|_\infty = \sup_{\mathbf{x} \in \mathbb{R}^p} \{|\Lambda_d^{\phi_n}(\mathbf{x}) - \Lambda_d^\phi(\mathbf{x})|\} \xrightarrow[n \rightarrow +\infty]{} 0 \quad \text{a.s.} \quad (1.52)$$

This theorem provides some convergence guarantees that justify the use of the empirical Christoffel function. In the context of data analysis, we are rarely aware of the probability measure that our data was generated from, which is why such guarantees are important.

Finally, a similar result as in Theorem 1.4.1. can be stated in terms of level sets. More precisely, we can certify that the level sets of the empirical and theoretical Christoffel function will asymptotically coincide, for any fixed degree. Once again, such a result puts us in a position to exploit the behavior of the Christoffel function when it comes to solving specific problems that arise in data analysis, which will be more detailed in the next part of the report.

# Chapter 2

## Application in machine learning

### 2.0.1 Introductory remarks

This section is going to be mainly focused on adapting the Christoffel function to the supervised learning context. Indeed, results from [1, 4, 11] have highlighted some remarkable capabilities of this function when it comes to recovering or approximating the graph of an unknown function. Since the classification task can be viewed as a specific task of the graph approximation, it becomes then natural to connect these ideas with the Christoffel function [5]. However, such an approach requires a certain amount of slight modifications to be made to the Christoffel function, which will be detailed in a sequel. This will allow us to define a classifier based on the Christoffel function and to evaluate its performance in various classification tasks.

### 2.1 Supervised learning

Suppose that we want to solve a classification problem of assigning the points in  $\mathbb{R}^p$  to one of the  $m \geq 2$  classes in the set  $\mathbf{Y} := \{1, \dots, m\}$ . Moreover, denote by  $\mathbf{X}_j \subset \mathbb{R}^p$  the set of all points that fall into the given class  $j \in \mathbf{Y}$ . Let  $\mathbf{X} = \bigcup_{j \in \mathbf{Y}} \mathbf{X}_j$  and suppose that  $\mathbf{X}$  is open with compact closure.

Let  $\mu$  be a joint probability distribution on the set  $\Omega = \mathbf{X} \times \mathbf{Y}$ . One may always disintegrate  $\mu$  into its marginal distribution  $\phi$  on  $\mathbf{X}$  and its conditional distribution  $\nu_{\mathbf{x}}$  on  $\mathbf{Y}$  given a point  $\mathbf{x} \in \mathbf{X}$ , so that  $d\mu(\mathbf{x}, y) = \nu_{\mathbf{x}}(dy)\phi(d\mathbf{x})$  [6].

If the points to classify are perfectly separable, i.e. they can not belong to more than one class, then we can assume that the marginal distribution assigns no mass to the intersection of the supports, meaning that  $\phi(\mathbf{X}_i \cap \mathbf{X}_j) = 0$ , for all  $(i, j) \in \mathbf{Y}^2$  with  $i \neq j$ . This allows us to write  $\mu = \sum_{j \in \mathbf{Y}} \mu_j$ , where for all  $j \in \mathbf{Y}$ , we have

$$d\mu_j(\mathbf{x}, y) = \delta_j(dy)\phi_j(d\mathbf{x}), \quad \text{with } \phi_j \in \mathcal{P}(\mathbf{X}_j). \quad (2.1)$$

This also implies that  $\phi = \sum_{j \in \mathbf{Y}} \phi_j$ , because  $\phi(A) = \mu(A \times \mathbf{Y}) = \sum_{j \in \mathbf{Y}} \mu_j(A \times \mathbf{Y}) = \sum_{j \in \mathbf{Y}} \phi_j(A)$ , for any Borel set  $A \in \mathcal{B}(\mathbf{X})$ .

Recall now that the goal of the supervised learning [5, 9, 12] is to recover a classifier  $f$  such that

$$f: \mathbf{X} \longrightarrow \mathbf{Y} \quad (2.2)$$

$$\mathbf{x} \longmapsto f(\mathbf{x}) := \begin{cases} \sum_{j=1}^m j \mathbf{1}_{\mathbf{X}_j}(\mathbf{x}), & \text{if } \mathbf{x} \in \mathbf{X} \\ 0, & \text{otherwise} \end{cases} \quad (2.3)$$

Thus, evaluating  $f$  at  $\mathbf{x} \in \mathbf{X}$  gives either the class of the point  $\mathbf{x}$  or zero (if the point belongs to different classes simultaneously).

As a byproduct of this construction, we obtain an alternative formulation of the underlying measure  $\mu$ , namely

$$d\mu(\mathbf{x}, y) = \delta_{f(\mathbf{x})}(dy)\phi(d\mathbf{x}) \quad (2.4)$$

Hence, the graph  $G := \{(\mathbf{x}, f(\mathbf{x})), \mathbf{x} \in \mathbf{X}\}$  of the classifier  $f$  is nothing else but the support of our underlying probability measure  $\mu$ . This is why we can identify the graph recovery problem with the problem of learning a classifier, and connect these two using the Christoffel function.

## 2.2 Christoffel-Darboux classifier

This part is quite technical but its purpose is to explain how the Christoffel function needs to be modified in order to encompass the problem of the supervised learning. Indeed, the usual moment matrix in the previously described classification framework will suffer from rank deficiency, requiring thus some additional refinements. Stated differently, the bi-linear form defined in (1.9) may not be positive definite, but we may still find an alternative way to circumvent the problem.

In this case, the support of the measure  $\mu$  defined in (2.4) is contained in a real algebraic variety of the space  $\mathbb{R}^p$ , which is why (1.9) fails to define a valid scalar product on  $L^2(\mu)$ . One way to observe this is by considering the polynomial  $Q \in \mathbb{R}[\mathbf{x}, y]$  such that  $Q(\mathbf{x}, y) := \prod_{i=1}^m (y - i)$ . Then,  $V := \{(\mathbf{x}, y) \in \mathbb{R}^p \times \mathbb{R} \mid Q(\mathbf{x}, y) = 0\} = \mathbb{R}^p \times \mathbf{Y}$  is the real algebraic set containing the support of the measure  $\mu$ . Thus, for any  $0 \neq P \in \mathbb{R}[\mathbf{x}, y]$  that vanishes on  $V$ , and represented by the vector of coefficients  $\mathbf{p}$ , we would have

$$\mathbf{p}^T \mathbf{M}_d(\phi) \mathbf{p} = \int_{\Omega} P^2 d\mu = 0, \quad (2.5)$$

justifying the non-invertibility of the moment matrix, so that equations (1.30) are not valid any more [4, 5]. In order to avoid these issues, one may define the Christoffel function as follows

$$\Lambda_d^\mu: V \longrightarrow \mathbb{R}_+ \quad (2.6)$$

$$(\mathbf{x}, y) \longmapsto \Lambda_d^\mu(\mathbf{x}, y) := \min_{P \in L_d^2(\mu)} \left\{ \int_{\Omega} P^2(\mathbf{z}) d\phi(\mathbf{z}), P(\mathbf{x}, y) = 1 \right\}, \quad (2.7)$$

where  $L_d^2(\mu) \subset \mathbb{R}[\mathbf{x}, y]$  is the space of polynomials on  $V$  of overall degree at most  $d$  (and at most  $m-1$  in the variable  $y$ ). A basis of such a set is thus given by  $\{\mathbf{x}^\alpha y^k \mid k \leq m-1, \sum_{i=1}^p \alpha_i + k \leq d\}$ . Then, the associated moment matrix is non-singular and we can retrieve the standard formulations of the Christoffel function on the set  $V$  [5].

Alternatively, we could perturb the measure  $\mu$  by considering  $\mu + \epsilon\mu_0$ , with  $\epsilon > 0$  and  $\mu_0$  a uniform probability measure on  $\mathbf{X} \times [0, m]$ . Then,

$$\Lambda_d^{\mu+\epsilon\mu_0}: \mathbb{R}^p \times \mathbb{R} \longrightarrow \mathbb{R}_+ \quad (2.8)$$

$$(\mathbf{x}, y) \longmapsto \Lambda_d^{\mu+\epsilon\mu_0}(\mathbf{x}, y) := \mathbf{v}_d(\mathbf{x}, y)^T M_d(\mu + \epsilon\mu_0)^{-1} \mathbf{v}_d(\mathbf{x}, y) \quad (2.9)$$

is well-defined and allows to approximate  $f$  via the following relationship [11]:

$$\hat{f}(\mathbf{x}) := \arg \min_y \Lambda_d^{\mu+\epsilon\mu_0}(\mathbf{x}, y). \quad (2.10)$$

However, a more intuitive and a more classification-oriented interpretation can be obtained by implementing another minor modification to the Christoffel function. Let us consider

$$\hat{\Lambda}_d^\mu: V \longrightarrow \mathbb{R}_+ \quad (2.11)$$

$$(\mathbf{x}, y) \longmapsto \hat{\Lambda}_d^\mu(\mathbf{x}, y) := \min_{P \in \mathcal{L}_d^2(\mu)} \left\{ \int_{\Omega} P^2(\mathbf{z}) d\phi(\mathbf{z}), P(\mathbf{x}, y) = 1 \right\}, \quad (2.12)$$

where  $\mathcal{L}_d^2(\mu) := \mathbb{R}_{d,m-1}[\mathbf{x}, y]$  is the space of all polynomials of degree at most  $d$  with respect to  $\mathbf{x}$  and at most  $m-1$  with respect to the variable  $y$ .

By construction,  $L_d^2(\mu) \subset \mathcal{L}_d^2(\mu) \subset L_{d+m-1}^2$ , so that for all  $(\mathbf{x}, y) \in V$ ,

$$\Lambda_{d+m-1}^\mu(\mathbf{x}, y) \leq \hat{\Lambda}_d^\mu(\mathbf{x}, y) \leq \Lambda_d^\mu(\mathbf{x}, y), \quad (2.13)$$

which is how these two constructions connect to each other. However, the advantage of the formulation in (2.12) lies in the fact that it possesses a closed form expression connected to the interpolation polynomials at the points  $\{1, \dots, m\}$ , which are defined through

$$y \mapsto \pi_j(y) := \frac{\prod_{i \neq j}(y-i)}{\prod_{i \neq j}(j-i)}, \quad j = 1, \dots, m \quad (2.14)$$

and satisfy  $\pi_j(i) = \delta_{ji}$ , for all  $i, j \in \mathbf{Y}$ .

This connection can be stated more precisely in the form of the following theorem

**Theorem 2.2.1.** [5] Let  $j \in \mathbf{Y}$  and let  $(P_\alpha^j)_{\alpha \in \mathbb{N}_d^p}$  be a family of polynomials in  $\mathbb{R}[\mathbf{x}]$  orthonormal with respect to  $\phi_j \in \mathcal{P}(\mathbf{X}_j)$ . Let  $\Lambda_d^{\phi_j}$  be the associated Christoffel function defined as in (1.28). Then,

1. The family  $(\pi_j P_\alpha^j)_{\alpha \in \mathbb{N}_d^p, j \in \mathbf{Y}}$  is an orthonormal basis of  $\mathcal{L}_d^2(\mu)$ ;
2.  $\hat{\Lambda}_d^\mu$  defined in (2.12) satisfies

$$(\hat{\Lambda}_d^\mu(\mathbf{x}, y))^{-1} = \sum_{j \in \mathbf{Y}} \pi_j(y)^2 \sum_{\alpha \in \mathbb{N}_d^p} P_\alpha^j(\mathbf{x})^2 = \sum_{j \in \mathbf{Y}} \delta_{y=j} (\Lambda_d^{\phi_j}(\mathbf{x}, y))^{-1}. \quad (2.15)$$

*Proof.*

1. The cardinality of the basis  $(\pi_j P_\alpha^j)_{\alpha \in \mathbb{N}_d^p, j \in \mathbf{Y}}$  is  $m \cdot s(d)$ , which corresponds to the dimension of the space  $\mathcal{L}_d^2(\mu) = \mathbb{R}_{d,m-1}[\mathbf{x}, y]$ . Using the fact that  $(\pi_j)_{j \in \mathbf{Y}}$  generates  $\mathbb{R}_{m-1}[y]$  and that  $(P_\alpha^j)_{\alpha \in \mathbb{N}_d^p}$  generates  $\mathbb{R}_d[\mathbf{x}]$ , we can easily deduce that every polynomial  $R \in \mathbb{R}_{d,m-1}[\mathbf{x}, y]$  can be written in as  $R(\mathbf{x}, y) = \sum_{\alpha \in \mathbb{N}_d^p, j \in \mathbf{Y}} p_j r_\alpha^j (\pi_j(y) P_\alpha^j(\mathbf{x}))$ . Finally, as long as  $i \neq j$ , we deduce from (2.14) that  $\pi_i(y) \pi_j(y) = 0$  inside the support of  $\mu$ , so that  $\int \pi_i(y) P_\alpha^i(\mathbf{x}) \pi_j(y) P_\beta^j(\mathbf{x}) d\mu(\mathbf{x}, y) = 0$  as well. On the other hand, if  $i = j$ , then

$$\int \pi_i(y)^2 P_\alpha^i(\mathbf{x}) P_\beta^i(\mathbf{x}) d\mu(\mathbf{x}, y) = \sum_{j \in \mathbf{Y}} \int \pi_i(y)^2 P_\alpha^i(\mathbf{x}) P_\beta^i(\mathbf{x}) d\mu_j(\mathbf{x}, y) \quad (\text{from (2.1)}) \quad (2.16)$$

$$= \int \pi_i(y)^2 P_\alpha^i(\mathbf{x}) P_\beta^i(\mathbf{x}) d\mu_i(\mathbf{x}, y) \quad (\text{by assumption}) \quad (2.17)$$

$$= \int P_\alpha^i(\mathbf{x}) P_\beta^i(\mathbf{x}) d\phi_i(\mathbf{x}, y) \quad (\text{from (2.1)}) \quad (2.18)$$

$$= \mathbf{1}_{\alpha=\beta} \quad (\text{by definition}) \quad (2.19)$$

which proves orthogonality.

2. Consider the monomial basis of  $\mathcal{L}_d^2(\mu)$ ,  $(\mathbf{x}^\alpha y^j)_{\alpha \in \mathbb{N}_d^p, 0 \leq j \leq m-1}$ , and denote by  $\hat{\mathbf{v}}(\mathbf{x}, y)$  the associated vector of monomials. Let  $\hat{M}_d(\mu)$  be the corresponding moment matrix defined as in (1.3). Following (1.29), we deduce that

$$\hat{\Lambda}_d(\mathbf{x}, y) = \min_{\mathbf{p} \in \mathbb{R}^{m \cdot s(d)}} \left\{ \mathbf{p}^T \hat{M}_d(\mu) \mathbf{p}, \mathbf{p}^T \hat{\mathbf{v}}_d(\mathbf{x}, y) = 1 \right\}. \quad (2.20)$$

We can then define the Lagrangian  $(\mathbf{p}, \lambda) \mapsto L(\mathbf{p}, \lambda) := \mathbf{p}^T \hat{M}_d(\mu) \mathbf{p} - \lambda(1 - \mathbf{p}^T \hat{\mathbf{v}}_d(\mathbf{x}, y))$ . The objective function being strictly convex, and the constraint being linear, we deduce that the first order condition for this optimization problem is both necessary and sufficient, and allows to characterize the optimal admissible solution  $(\mathbf{p}^*, \lambda^*)$  as follows

$$2\hat{M}_d(\mu)\mathbf{p}^* = \lambda^* \hat{\mathbf{v}}_d(\mathbf{x}, y), \quad (2.21)$$

which after multiplying by  $\mathbf{p}^*$  yields

$$\lambda^* = 2\hat{\Lambda}_d(\mathbf{x}, y) \quad \text{and} \quad \mathbf{p}^* = \hat{\Lambda}_d(\mathbf{x}, y)(\hat{M}_d(\mu))^{-1} \hat{\mathbf{v}}_d(\mathbf{x}, y). \quad (2.22)$$

This implies that the optimal polynomial in the definition of  $\hat{\Lambda}_d(\mathbf{x}, y)$  writes, for all  $(\mathbf{u}, z) \in \mathbb{R}^p \times \mathbf{Y}$ ,

$$\begin{aligned} p^*(\mathbf{u}, z) &= \hat{\mathbf{v}}_d(\mathbf{x}, y)^T \hat{\Lambda}_d(\mathbf{x}, y) (\hat{M}_d(\mu))^{-1} \hat{\mathbf{v}}_d(\mathbf{x}, y) \\ &= \hat{\Lambda}_d(\mathbf{x}, y) \sum_{\alpha \in \mathbb{N}_d^p, j \in \mathbf{Y}} \pi_j(y) \pi_j(z) P_\alpha^j(\mathbf{x}) P_\alpha^j(\mathbf{u}) \quad (\text{from (1.11) and Proposition 1.4.1}) \end{aligned} \quad (2.23)$$

By evaluating  $p^*$  at  $(\mathbf{x}, y)$ , we deduce that

$$p^*(\mathbf{x}, y) = 1 = \hat{\Lambda}_d(\mathbf{x}, y) \sum_{\alpha \in \mathbb{N}_d^p, j \in \mathbf{Y}} \pi_j(y)^2 P_\alpha^j(\mathbf{x})^2 \quad (2.24)$$

$$= \hat{\Lambda}_d(\mathbf{x}, y) \sum_{j \in \mathbf{Y}} \pi_j(y)^2 \sum_{\alpha \in \mathbb{N}_d^p} P_\alpha^j(\mathbf{x})^2 \quad (2.25)$$

$$= \hat{\Lambda}_d(\mathbf{x}, y) \sum_{j \in \mathbf{Y}} \pi_j(y)^2 (\Lambda_d^{\phi_j}(\mathbf{x}))^{-1} \quad (\text{use definition}) \quad (2.26)$$

$$= \hat{\Lambda}_d(\mathbf{x}, y) \sum_{j \in \mathbf{Y}} \delta_{y=j} (\Lambda_d^{\phi_j}(\mathbf{x}))^{-1} \quad (\text{from (2.14)}) \quad (2.27)$$

which are exactly those equalities stated in (2.15).  $\square$

Thus, if  $y \in \mathbf{Y}$ , we can easily and simply express the modified Christoffel function in terms of the class-specific standard Christoffel functions  $\Lambda_d^{\phi_j}$ , for  $j = 1, \dots, m$ .

## 2.2.1 Consistency result

The expression that was derived in Theorem 2.2.1. is what motivates the following definition.

**Definition 2.2.1.** The *Christoffel-Darboux classifier* (CDC) of order  $d \geq 1$  is a function  $\hat{f}_d$  such that

$$\forall \mathbf{x} \in \mathbf{X}, \quad \hat{f}_d(\mathbf{x}) := \arg \max_{j \in \llbracket 1, m \rrbracket} \Lambda_d^{\phi_j}(\mathbf{x}). \quad (2.28)$$

That this definition is natural follows from the following observation. Fix  $k \in \mathbf{Y}$  and  $\mathbf{x} \in \mathbf{X}_k$ . From (1.37) and (1.47), we know that, if we increase the order  $d$ , then  $\Lambda_d^{\phi_j}(\mathbf{x})$  should decrease exponentially fast towards zero for all  $j \neq k$ . On the other hand,  $\Lambda_d^{\phi_k}(\mathbf{x})$  should manifest a decay which is not faster than polynomial.

So, if the degree is sufficiently large (say higher than a threshold  $d_x \in \mathbb{N}$ ), then after some time, we should necessarily have  $\Lambda_d^{\phi_j}(\mathbf{x}) < \Lambda_d^{\phi_k}(\mathbf{x})$ , for all  $d \geq d_x$  and all  $j \neq k$ . Consequently, those class-specific Christoffel functions can be interpreted as score-determining, so that every new point is classified to the class that yields the largest class-specific Christoffel function.

That this procedure represents a consistent way of defining a classifier is stated in the following theorem:

**Theorem 2.2.2.** [5] Let  $j \in \llbracket 1, m \rrbracket$ ,  $\epsilon > 0$  and  $\mathbf{X}_j^\epsilon := \{\mathbf{x} \in \mathbf{X}_j \mid d(\mathbf{x}, \partial \mathbf{X}_j) > \epsilon\}$ . Assume that for all  $j \in \llbracket 1, m \rrbracket$ ,  $\phi_j$  is absolutely continuous with respect to the restriction of the Lebesgue measure on  $\mathbf{X}_j$ . Then, there exists  $d_{x,\epsilon} \in \mathbb{N}$  such that for all  $d \geq d_{x,\epsilon}$

$$\hat{f}_d(\mathbf{x}) = j, \quad \forall \mathbf{x} \in \mathbf{X}_j^\epsilon. \quad (2.29)$$

Due to the high number of technicalities, we choose to omit the proof of the Theorem 2.2.2. Nonetheless, the main message that should be taken from this theorem is that, by adjusting the order of the classifier  $\hat{f}_d$ , one can make sure that all the points, which are sufficiently far away from the class boundary, get correctly classified, which is quite a desirable feature.

It is important to notice that the classifier we defined in (2.28) remains purely theoretical, since we rarely know those population probability distributions  $(\phi_k)_{k \in \mathbf{Y}}$ . What we have instead at our disposal is the finite simple of labeled data  $\{(\mathbf{x}_i, y_i) \in \mathbf{X} \times \mathbf{Y}\}_{i=1}^N$  often referred to as the *training data set*. These data points can be then used to construct empirical counterparts of the quantities defined in (2.28).

Indeed, we may consider the sequence of empirical measures  $(\phi_{k,N})_{k \in \mathbf{Y}}$  and their associated empirical Christoffel functions  $(\Lambda_d^{\phi_{k,N}})_{k \in \mathbf{Y}}$ . The natural question that this construction raises is whether these discretizations still provide good approximations of the support of the underlying probability measures.

To answer this question, suppose that the training sets consists of  $N$  data points from each of the  $m$  classes. Compute for all  $k \in \mathbf{Y}$ ,  $\phi_{k,N} := \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^k$ , and let  $(\Lambda_d^{\phi_{k,N}}(\mathbf{x}))^{-1} = \mathbf{v}_d(\mathbf{x})^T (M_d(\phi_{k,N}))^{-1} \mathbf{v}_d(\mathbf{x})$ . These functions are indeed well-defined as soon as the sample size is large enough and  $d$  is chosen carefully so that the empirical moment matrix is indeed invertible [2].

Then, the empirical analogue of the CDC is given by

$$\forall \mathbf{x} \in \mathbf{X}, \quad \hat{f}_d^N(\mathbf{x}) := \arg \max_{j \in \llbracket 1, m \rrbracket} \Lambda_d^{\phi_{j,N}}(\mathbf{x}). \quad (2.30)$$

We have the following result:

**Theorem 2.2.3.** [5] Let  $j \in \llbracket 1, m \rrbracket$ ,  $\epsilon > 0$  and  $\mathbf{X}_j^\epsilon := \{\mathbf{x} \in \mathbf{X}_j \mid d(\mathbf{x}, \partial \mathbf{X}_j) > \epsilon\}$ . Let  $(\mathbf{x}_i^j)_{i=1}^N \subset \mathbf{X}_j$  be an i.i.d. sample drawn from an absolutely continuous distribution  $\phi_j \in \mathcal{P}(\mathbf{X}_j)$ . Then, there exists  $d_{x,\epsilon} \in \mathbb{N}$  such that for all  $d \geq d_{x,\epsilon}$ , and for sufficiently large sample size  $N$ ,

$$\hat{f}_d^N(\mathbf{x}) = j, \quad \forall \mathbf{x} \in \mathbf{X}_j^\epsilon, \quad (2.31)$$

almost surely with respect to the random samples.

*Proof.* Let  $j \in \llbracket 1, m \rrbracket$  and  $\epsilon > 0$ . Let  $\mathbf{x} \in \mathbf{X}_j^\epsilon$ . From the Theorem 2.2.2, we deduce that it is possible to find a real constant  $c_j > 0$  and  $d_{\mathbf{x},\epsilon} \in \mathbb{N}$  such that  $\Lambda_d^{\phi_j}(\mathbf{x}) - \Lambda_d^{\phi_k}(\mathbf{x}) > c_j$ , for all  $k \neq j$  and all  $d \geq d_{\mathbf{x},\epsilon}$ .

Furthermore, using the Theorem 1.4.1, we know that, for all  $k \in \llbracket 1, m \rrbracket$  and all  $d \in \mathbb{N}^*$ ,

$$\sup_{\mathbf{x} \in \mathbb{R}^p} \{|\Lambda_d^{\phi_{k,N}}(\mathbf{x}) - \Lambda_d^{\phi_k}(\mathbf{x})|\} \xrightarrow{N \rightarrow +\infty} 0 \quad (2.32)$$

almost surely with respect to the random samples  $(\mathbf{x}_i^k)_{i=1}^N \subset \mathbf{X}_k$ . If we combine these two remarks, we deduce that for all  $d > d_{\mathbf{x},\epsilon}$  and all  $k \neq j$ ,

$$\Lambda_d^{\phi_{j,N}}(\mathbf{x}) > \Lambda_d^{\phi_{k,N}}(\mathbf{x}) + c_j, \quad (2.33)$$

provided that  $N$  is sufficiently large, which is the same as saying that  $\hat{f}_d^N(\mathbf{x}) = j$ .  $\square$

In other words, Theorem 2.2.3 provides some theoretical guarantees that could be used to justify the use of the CDC in the standard supervised learning framework. Implementing and evaluating this classifier is what the next sections are going to be mostly based on.

## 2.3 Empirical evaluations

At this stage, we are interested in implementing numerically the CDC classifier and measuring its performance on various data sets. More precisely, the results that will be presented below are obtained thanks to our preliminary code written in Python.

A significant amount of preliminary computations was necessary in order to successfully implement the CDC classifier. For example, the following functions were created:

- *gen\_basis(p,d)* – generates the  $p$ -variate monomial basis of order  $d$  w.r.t. the graded lexicographic order;
- *MoEM(data,d)* – computes the empirical moment matrix (indexed by monomials) of order  $d$  associated to the points in *data*;
- *Christoffel(point,d,data)* – evaluates at *point* the empirical Christoffel function of order  $d$  associated to the cloud of points in *data*.

Next, inspired by the already available classifiers from the *scikit-learn* library, we create a new **class CDKC()** of objects which will be responsible for implementing the formula in (2.30).

Some of the **main attributes** associated to this class of objects are:

- *emommatrix* - computes the empirical moment matrix from the input data;
- *predict(tobeclassified)* - assigns labels to the points in *tobeclassified*;
- *score(tobeclassified,true\_label)* - percentage of correctly classified points;
- *confusion\_matrix(true\_label,prediction)* - confusion matrix computation.

It is important to mention that the code has not been optimized, but rather used in order to obtain some benchmark comparisons. Furthermore, we emphasize the fact that the CDC implementation is dependant on only one tuning parameter, namely the degree  $d$  of the involved Christoffel functions.

### 2.3.1 *Iris* data

We first evaluate the performance of the CDC classifier on the very famous *Iris data* set. This data set contains four features (length and width of sepals and petals) of 50 instances of three species of *Iris* flower, i.e.  $\mathbf{Y} = \{1, 2, 3\}$  and  $N = 50$ , so that we have 150 labeled data points that were split into training and test categories via standard procedures.

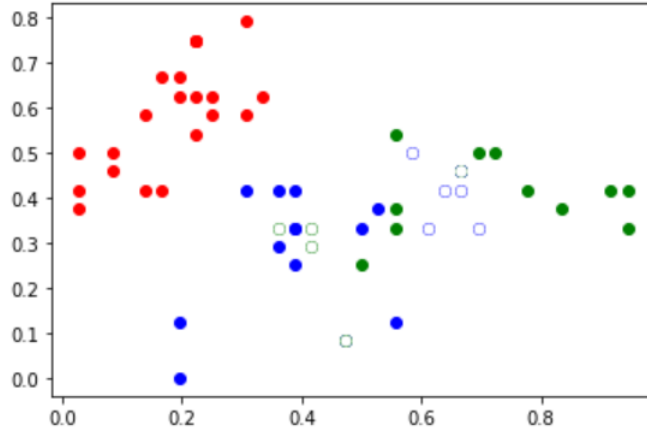


Figure 2.1: Scaled *Iris* data - CDC classification example with  $d = 4$  ; points whose interior is not coloured are those that were miss-classified.

What we can observe from Figure 2.1. is that the red points were all correctly classified. Indeed, miss-classification occurred among blue and green points only since their supports are not disjoint, which violates the assumptions that we used in our theoretical framework. These intersecting supports make the exact classification impossible. Notice also that only two out of four features were used to solve the classification task.

Another interesting observation, which is rather intuitive, is that by considering all the features, one may increase the percentage of the correctly classified data. Indeed, in the 4-dimensional space, this data becomes much more easily separable.

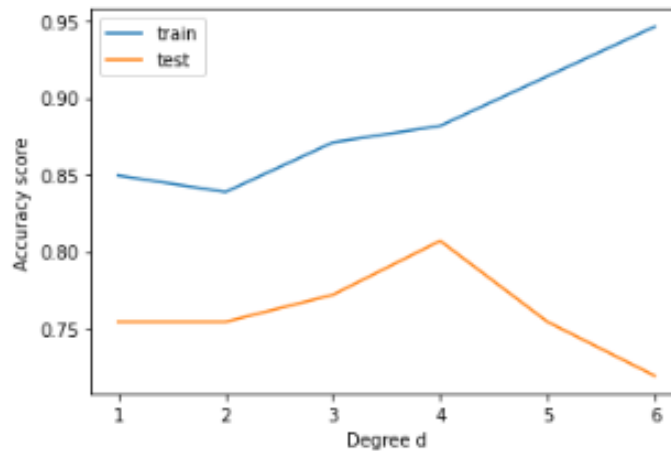


Figure 2.2: Score of the CDC classifier on the test and training sets as a function of the tuning parameter  $d$ .



As expected, if one increases the degree of the CDC classifier, its performance also increases and stabilizes around a certain threshold, if the data is not perfectly separable. However, by doing so, it seems that the levels sets of the Christoffel functions involved in the CDC expression become too specific to the data from the training set, and thus fail to generalize to the data in the test set. This over-fitting tendency seems to be captured by the Figure 2.2. However, an alternative explanation would be that the rather small sample size and degree are not sufficient to make the moment matrix invertible, so that this behavior is a numerical artefact resulting from the inversion of ill-conditioned matrices (the smallest eigenvalues of label-specific moment matrices are of order  $10^{-10}$ ).

Let us now compare the CDC classifier with other commonly used classifiers:

Comparison for $p = 2$ features		
Classifier	Precision (%)	Time (ms)
CDC ( $d = 4$ )	81	21.8
kNN	79	2.31
LDA	75	2.54
SVM	77	3.19
Tree	74	2.67

Table 2.1: Mean execution time and accuracy comparison

First, we notice that, even for small values of the tuning parameter, the CDC classifier is able to achieve great performance compared to the other commonly used classifiers. Furthermore, we observe that this method is much more costly to implement numerically, since the mean execution time is higher than what many other classifiers require. Moreover, by considering all the four features, the CDC classifier is able to correctly separate all the points by using only the moments up to order four, i.e.,  $d = 2$ . However, the method’s computational speed deteriorates significantly (200 ms on average), while the computing time of other classifiers remains quite stable. This indicates that the current implementation of the CDC classifier suffers from important dimensionality issues, so that the classifier is suitable for small-size problems only.

### 2.3.2 Toy data sets

We study the behaviour of the CDC classifier on various artificial data sets.

We start by considering the famous generic *Two moons* example. We control the level of the noise so that the two moon-like shapes remain disjoint (separable). Some of the obtained results are summarized by the means of the following figure:

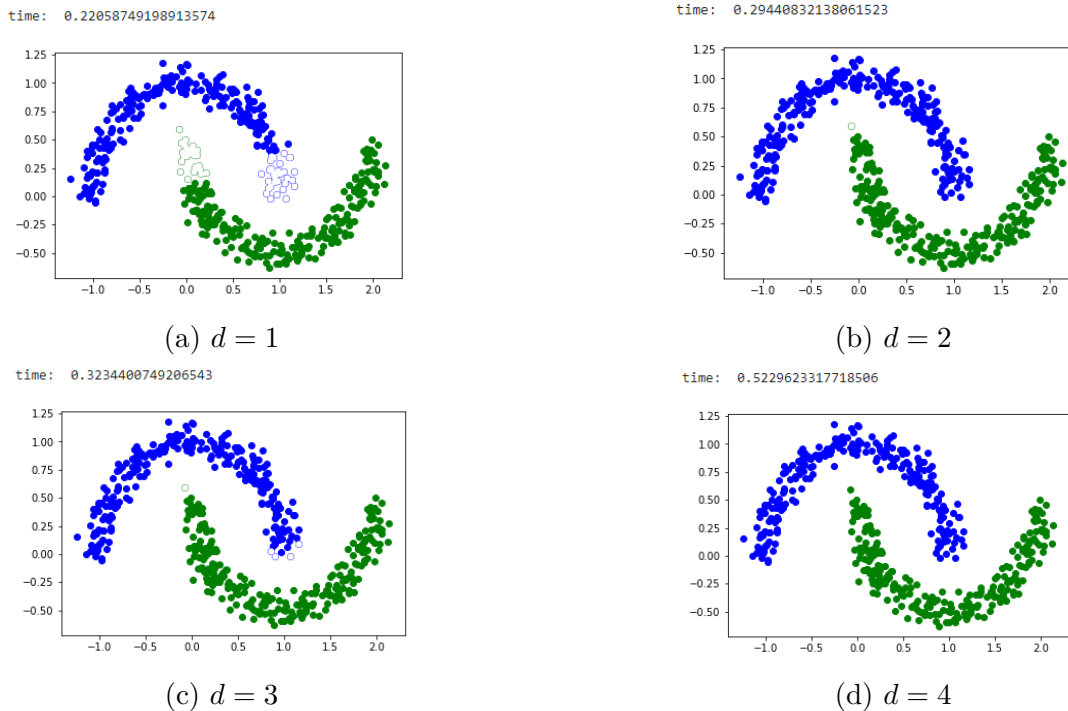


Figure 2.3: Performance of the CDC classifier on the *Two moons* data set, with  $p = 2$  and  $N = 2000$ ; Computation time is expressed in seconds.

Figure 2.3. displays some convergence properties of the CDC classifier that are in line with Theorem 2.2.3. Indeed, it suffices to consider the Christoffel functions of order not higher than four (i.e. empirical moments of order eight) in order to extract the level sets that perfectly separate the points in the depicted data cloud.

On the same data set, the LDA classifier was much faster but achieved only 88% accuracy, demonstrating some advantages of the CDC classifier in the context of non-linearly separable data sets. SVM classifier, on the other hand, manages to achieve the absolute efficiency by consuming somehow comparable amount of time ( $123ms$ ) as  $\hat{f}_d^N$ .

Similar results are obtained for other data shapes.

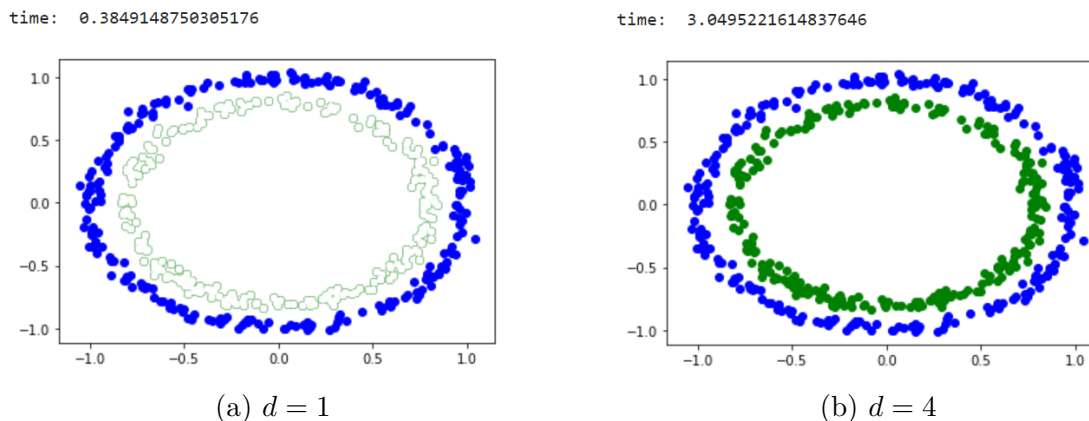


Figure 2.4: Embedded disks and CDC classifier,  $p = 2$  and  $N = 2000$ .

What Figure 2.4. displays is a fifty percent accuracy gain obtained by considering polynomials of degree 8 instead of the much simpler degree two polynomials.

Another observation that can be made at this stage of analysis is related to the dependency of the CDC classifier on the sample size. It turns out that larger  $N$  is much less harmful in terms of computational complexity than, for example, larger  $p$ .

Let us end this section by including some multi-class performance evaluations.

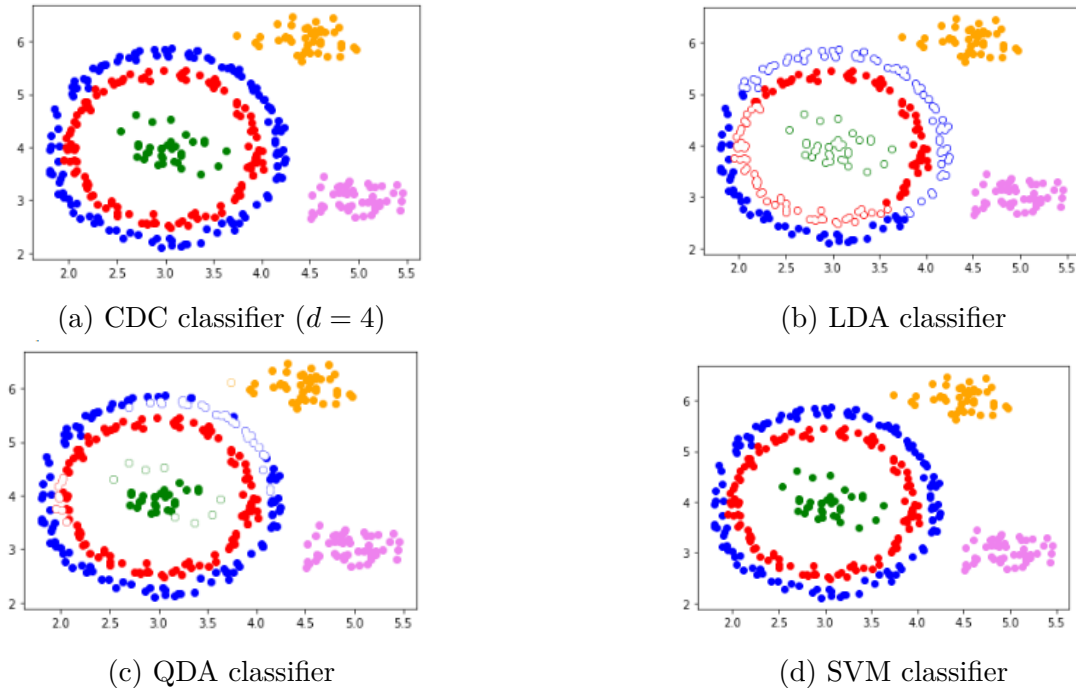


Figure 2.5: Comparison of various classifiers,  $\mathbf{Y} = \llbracket 1, 5 \rrbracket, p = 2$ .

Even in the case of non-binary classification, the Christoffel-Darboux classifier seems to perform pretty well in terms of accuracy.

With only one tuning parameter, which has a very straightforward interpretation, the CDC classifier provides us with a classification tool that is deeply connected to the data generating process. This connection is established through the support identification, which takes place when the Christoffel-Darboux classifier is computed. This nice interpretability and quite natural parameter dependency are the main advantages of this classifier. On the other hand, being extremely computationally expensive in higher dimensions is the main drawback of this classification method.

## 2.4 Extensions

As emphasized previously, there are still many improvements to be made to the proposed classification method. Circumventing dimensionality issues and alleviating computational difficulties are among the main directions for the future work.

For example, it would be interesting to implement the CDC classifier using the relationship established in (1.29). This alternative way of computing the Christoffel function would make it

unnecessary to invert moment matrices of high dimensions, manifesting probably better numerical behaviour. However, this new method requires solving a completely new convex quadratic programming problem (which could be done efficiently using many solvers) whenever one needs to evaluate the Christoffel function at a different point. What is the trade-off between these two separate techniques remains to be investigated.

Alternatively, one may try to combine (1.29) or the variational formulation in (1.28) with neural networks and learn the Christoffel function. Indeed, in (1.28) one solves an optimization problem over the set of polynomials of a given degree. As such, that formulation is not tractable. But, one could consider optimizing the same objective function over the set of neural networks. Such an approach would result in an approximate CDC classifier which would not suffer from the curse of dimensionality. Since polynomials are continuous functions, and neural networks possess very interesting expressivity properties, one may hope to recover an optimal solution with desired error margin.

Additional research directions would include generalizing this approach to other tasks in the realm of data analysis. For example, it would be natural to examine how the Christoffel-function-based-approach behaves in the unsupervised learning framework (clustering). Here, one does not know *a priori* how many clusters actually exist, i.e. labels are not available. An efficient clustering algorithm should thus be able to deduce from the entire cloud of points (delineated by the level sets of one polynomial) how many natural sub-clouds exist, i.e., to decompose one *big* Christoffel function into multiple *small* Christoffel functions [6, 14]. Measuring the quality of such decompositions, while keeping the number of tuning parameters as small as possible is one of the main challenges.

# Conclusion

In this report, we have provided some results related to the use of the Christoffel function in machine learning. More precisely, we have investigated how this tool from approximation theory domain can be adapted to handle supervised learning, i.e. classification tasks. Some of those results were purely theoretical, like in the first part, but found their use in practical implementations described in the second part.

As already mentioned, the scientific community has been aware for a very long time of many remarkable properties of the Christoffel-Darboux kernel, that were afterwards exploited mainly in the study of orthogonal polynomials. Recently, it was shown that this tool can be used in statistics as well, since the level sets of the Christoffel functions associated to empirical measures identify quite well the supports of those measures. This rather simple statement has motivated the content of the first part of the report. Indeed, getting familiar with the properties of the Christoffel function is a prerequisite for properly manipulating those properties when solving data science problems.

Furthermore, we have focused on solving the classification problem with a novel classifier based on the Christoffel function.

It turns out that a slightly modified Christoffel function represents the cornerstone for developing a classifier with some interesting properties. Indeed, we have demonstrated that such a classifier is quite easy to interpret, since it is profoundly connected to the data generating process. Moreover, it depends on only one tuning parameter, which is also a very much desired feature. In addition, some consistency-like results were stated that justify this approach from the theoretical point of view.

Finally, we have discussed numerical implementations and measured performance qualities of the proposed estimator. We have seen that, in terms of accuracy, the CDC classifier competes very well with many other state-of-the-art classifiers, and on different data sets. On the other hand, this classifier involves some pretty heavy computations that fail to make it suitable for the real-life classification problems, which are more than often high-dimensional. In any case, this approach represents a very promising direction for future works that could be useful in many machine learning applications.

# Bibliography

- [1] D. Henrion, J. B. Lasserre (2022) *Graph Recovery from Incomplete Moment Information*. Constructive Approximations <https://doi.org/10.1007/s00365-022-09563-8>
- [2] E. Pauwels, M. Putinar, J. B. Lasserre (2018). *Data analysis from empirical moments and the Christoffel function*. arXiv:1810.08480 [stat.ML]
- [3] E. Pauwels, J. B. Lasserre (2016). *Sorting out typicality with the inverse moment matrix SOS polynomial*. Advances in Neural Information Processing Systems 29 (NIPS 2016)
- [4] J. B. Lasserre, E. Pauwels, M. Putinar (2022). *The Christoffel–Darboux Kernel for Data Analysis*. Cambridge University Press
- [5] J. B. Lasserre (2022). *On the Christoffel function and classification in data analysis*. arXiv:2203.14571 [math.OC]
- [6] J. B. Lasserre (2022). *A disintegration of the Christoffel function*. hal03624003
- [7] J. B. Lasserre (2001). *Global optimization with polynomials and the problem of moments*. SIAM Journal on optimization, 11(3), 796-817.
- [8] N. Aronszajn (1950). *Theory of Reproducing Kernels*. Transactions of the American Mathematical Society, Vol. 68, No. 3 (May, 1950), pp. 337-404
- [9] P. Neval (1986). *Géza Freud, orthogonal polynomials and Christoffel functions. A case study*. Journal of Approximation Theory
- [10] S. B. Kotsiantis, I. Zaharakis, and P. Pintelas (2007). *Supervised machine learning: A review of classification techniques*. Emerging artificial intelligence applications in computer engineering 160, no. 1: 3-24.
- [11] S. Marx, E. Pauwels, T. Weisser, D. Henrion and J. B. Lasserre (2021). *Semi-algebraic Approximation Using Christoffel–Darboux Kernel*. <https://doi.org/10.1007/s00365-021-09535-4>
- [12] S. L. Brunton and J. N. Kutz (2019). *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, Cambridge, UK
- [13] S. T. Belinschi, V. Magron and V. Vinnikov (2021) . *Non-commutative Christoffel-Darboux Kernels*. arXiv preprint, arXiv:2106.06212.
- [14] J. W. Helton, J. B. Lasserre and M. Putinar (2008). *Measures with zeros in the inverse of their moment matrix*. The Annals of Probability, 36(4), 1453-1471.