

Sorting out typicality with the inverse moment matrix SOS polynomial

Authored by:

Edouard Pauwels
Jean B. Lasserre

Abstract

We study a surprising phenomenon related to the representation of a cloud of data points using polynomials. We start with the previously unnoticed empirical observation that, given a collection (a cloud) of data points, the sublevel sets of a certain distinguished polynomial capture the shape of the cloud very accurately. This distinguished polynomial is a sum-of-squares (SOS) derived in a simple manner from the inverse of the empirical moment matrix. In fact, this SOS polynomial is directly related to orthogonal polynomials and the Christoffel function. This allows to generalize and interpret extremality properties of orthogonal polynomials and to provide a mathematical rationale for the observed phenomenon. Among diverse potential applications, we illustrate the relevance of our results on a network intrusion detection task for which we obtain performances similar to existing dedicated methods reported in the literature.

1 Paper Body

Capturing and summarizing the global shape of a cloud of points is at the heart of many data processing applications such as novelty detection, outlier detection as well as related unsupervised learning tasks such as clustering and density estimation. One of the main difficulties is to account for potentially complicated shapes in multidimensional spaces, or equivalently to account for non standard dependence relations between variables. Such relations become critical in applications, for example in fraud detection where a fraudulent action may be the dishonest combination of several actions, each of them being reasonable when considered on their own. Accounting for complicated shapes is also related to computational geometry and nonlinear algebra applications, for example integral computation [11] and reconstruction of sets from moments data [6, 7, 12]. Some of these problems have connections and potential applications in machine learning. The work presented in this paper brings together ideas from both disciplines, leading to a method which allows to encode in a simple manner the global shape and spatial concentration of points within a cloud. We start

with a surprising (and apparently unnoticed) empirical observation. Given a collection of points, one may build up a distinguished sum-of-squares (SOS) polynomial whose coefficients (or Gram matrix) is the inverse of the empirical moment matrix (see Section 3). Its degree depends on how many moments are considered, a choice left to the user. Remarkably its sublevel sets capture much of the global shape of the cloud as illustrated in Figure 3. This phenomenon is not incidental as illustrated in many additional examples in Appendix A. To the best of our knowledge, this observation has remained unnoticed and the purpose of this paper is to report this empirical finding to the machine learning community and provide first elements toward a mathematical understanding as well as potential machine learning applications. 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

146370 67380
30670
30 67 67 38 0 0 14 63 70
35
13810 6860
3570
1850
950 460
340
210
333
80 10
13
70 18 50 95 81 0 0 68 60
333340
30 15
138
10
357 0 6860
185
0
950

Figure 1: Left: 1000 points in \mathbb{R}^2 and the levelsets of the corresponding inverse moment matrix SOS polynomial $Q_{p,d}$ ($d = 4$). The level set $p+d$, which corresponds to the average value of $Q_{p,d}$, is represented in red. Right: 1040 points in \mathbb{R}^2 with size and color proportional to the value of inverse moment matrix SOS polynomial $Q_{p,d}$ ($d = 8$). The proposed method is based on the computation of the coefficients of a very specific polynomial which depends solely on the empirical moments associated with the data points. From a practical perspective, this can be done via a single pass through the data, or even in an online fashion via a sequence of efficient Woodbury updates. Furthermore the computational cost of evaluating the polynomial does not depend on the number of data points which is a crucial difference with existing non-parametric methods such as nearest neighbors or kernel based methods [3]. On

the other hand, this computation requires the inversion of a matrix whose size depends on the dimension of the problem (see Section 3). Therefore, the proposed framework is suited for moderate dimensions and potentially very large number of observations. In Section 4 we first describe an affine invariance result which suggests that the distinguished SOS polynomial captures very intrinsic properties of clouds of points. In a second step, we provide a mathematical interpretation that supports our empirical findings based on connections with orthogonal polynomials [5]. We propose a generalization of a well known extremality result for orthogonal univariate polynomials on the real line (or the complex plane) [16, Theorem 3.1.2]. As a consequence, the distinguished SOS polynomial of interest in this paper is understood as the unique optimal solution of a convex optimization problem: minimizing an average value over a structured set of positive polynomials. In addition, we revisit [16, Theorem 3.5.6] about the Christoffel function. The mathematics behind provide a simple and intuitive explanation for the phenomenon that we empirically observed. Finally, in Section 5 we perform numerical experiments on KDD cup network intrusion dataset [13]. Evaluation of the distinguished SOS polynomial provides a score that we use as a measure of outlyingness to detect network intrusions (assuming that they correspond to outlier observations). We refer the reader to [3] for a discussion of available methods for this task. For the sake of a fair comparison we have reproduced the experiments performed in [18] for the same dataset. We report results similar to (and sometimes better than) those described in [18] which suggests that the method is comparable to other dedicated approaches for network intrusion detection, including robust estimation and Mahalanobis distance [8, 10], mixture models [14] and recurrent neural networks [18].

2

Multivariate polynomials, moments and sums of squares

Notations: We fix the ambient dimension to be p throughout the text. For example, we will manipulate vectors in \mathbb{R}^p as well as p -variate polynomials with real coefficients. We denote by X a set of p variables X_1, \dots, X_p which we will use in mathematical expressions defining polynomials. We identify monomials from the canonical basis of p -variate polynomials with their exponents in \mathbb{N}^p : we associate P_p to $\alpha = (\alpha_i)_{i=1 \dots p} \in \mathbb{N}^p$ the monomial $X := X_1^{\alpha_1} X_2^{\alpha_2} \dots X_p^{\alpha_p}$ which degree is $\deg(\alpha) := \sum_{i=1}^p \alpha_i$. We use the expressions $|\alpha|$ and $|\alpha|$ to denote the graded lexicographic order, a well ordering over p -variate monomials. This amounts to, first, use the canonical order on the 2

degree and, second, break ties in monomials with the same degree using the lexicographic order with $X_1 = a, X_2 = b, \dots$. For example, the monomials in two variables X_1, X_2 , of degree less or equal to 3 listed in this order are given by: 1, X_1 , X_2 , X_1^2 , $X_1 X_2$, X_2^2 , X_1^3 , $X_1^2 X_2$, $X_1 X_2^2$, X_2^3 . We denote by N_p , the set $\{\alpha \in \mathbb{N}^p; \deg(\alpha) \leq d\}$ ordered by $|\alpha|$. $R[X]$ denotes the set of p -variate polynomials: linear combinations of monomials with real coefficients. The degree of a polynomial is the highest of the degrees of its monomials with nonzero coefficients¹. We use the same notation, $\deg(\alpha)$, to denote the degree of a polynomial or of an element of N_p . For $d \in \mathbb{N}$, $R_d[X]$ denotes the set of p -variate polynomials of degree less or equal to d . We set $s(d) = p + d$, the

number of monomials of degree less or equal to d . We will denote by $\text{vd}(X)$ the vector of monomials of degree less or equal to d sorted by ?gl . We let $\text{vd}(X) := (X^{\text{?}})^{\text{?Np}} \text{? Rd}[X]_{\text{s}(d)}$. With this notation, d we can write a polynomial $P \text{? Rd}[X]$ as follows $P(X) = \text{hp, vd}(X)^{\text{i}}$ for some real vector of coefficients $p = (p^{\text{?}})^{\text{?Np}} \text{? Rs}(d)$ ordered using ?gl . Given $x = (x_i)_{i=1\dots p} \text{? Rp}$, $P(x)$ denotes d the evaluation of P with the assignments $X_1 = x_1, X_2 = x_2, \dots, X_p = x_p$. Given a Borel probability measure ? and $\text{?} \text{? Np}$, $y^{\text{?}}(\text{?})$ denotes the moment ? of ? : $y^{\text{?}}(\text{?}) = \text{Rp } x^{\text{?}} d^{\text{?}}(x)$. Throughout the paper, we will only consider measures of which all moments are finite. Moment matrix: Given a Borel probability measure ? on Rp , the moment matrix of ? , $M_d(\text{?})$, is a matrix indexed by monomials of degree at most d ordered by ?gl . For $\text{?}, \text{?} \text{? Npd}$, the corresponding entry in $M_d(\text{?})$ is defined by $M_d(\text{?})^{\text{?},\text{?}} := y^{\text{?}} + y^{\text{?}}(\text{?})$, the moment $\text{?} + \text{?}$ of ? . When $p = 2$, letting $y^{\text{?}} = y^{\text{?}}(\text{?})$ for $\text{?} \text{? N24}$, we have

$$\begin{aligned} M_2(\text{?}) : \\ & \begin{matrix} 1 & X_1 & X_2 & X_{12} & X_1 & X_2 & X_{22} \\ 1 & & & & & & \\ X_1 & & & & & & \\ X_2 & & & & & & \\ X_{12} & & & & & & \\ X_1 & X_2 & & & & & \\ 1 & y_{10} & y_{01} & y_{20} & y_{11} & y_{02} & \\ y_{10} & y_{20} & y_{11} & y_{30} & y_{21} & y_{12} & \\ y_{01} & y_{11} & y_{02} & y_{21} & y_{12} & y_{03} & \\ y_{20} & y_{30} & y_{21} & y_{40} & y_{31} & y_{22} & \\ y_{11} & y_{21} & y_{12} & y_{31} & y_{22} & y_{13} & \\ X_{22} & y_{02} & y_{12} & y_{03} & y_{22} & y_{13} & y_{04} \end{matrix} \end{aligned}$$

$M_d(\text{?})$ is positive semidefinite for all $d \text{? N}$. Indeed, for any $\text{Rp } \text{? Rs}(d)$, let $P \text{? Rd}[X]$ be the polynomial with vector of coefficients p , we have $p^{\text{T}} M_d(\text{?}) p = \text{Rp } P^2(x) d^{\text{?}}(x) \text{? 0}$. Furthermore, R we have the identity $M_d(\text{?}) = \text{Rp } \text{vd}(x) \text{vd}(x)^{\text{T}} d^{\text{?}}(x)$ where the integral is understood elementwise. Sum of squares (SOS): We denote by $\text{?[X]} \text{? R}[X]$ (resp. $\text{?d}[X] \text{? Rd}[X]$), the set of polynomials (resp. polynomials of degree at most d) which can be written as a sum of squares of polynomials. Let $P \text{? R}^{2m}[X]$ for some $m \text{? N}$, then P belongs to $\text{?}^{2m}[X]$ if there exists a finite $J \text{? N}$ and a family of polynomials $P_j \text{? R}^m[X]$, $j \text{? J}$, such that $P = \sum_{j=1}^J P_j^2$. It is obvious that sum of squares polynomials are always nonnegative. A further interesting property is that this class of polynomials is connected with positive semidefiniteness. Indeed, P belongs to $\text{?}^{2m}[X]$ if and only if $\text{?Q } \text{? Rs}(m) \text{?s}(m)$, $Q \geq 0$, $P(x) = \text{vd}(x)^{\text{T}} Q \text{vd}(x)$, $\text{?x } \text{? Rp}$. As a consequence, every positive semidefinite matrix $Q \text{? R } \text{?}^{2m}[X]$ by using the representation in (1).

$$\begin{aligned} & \begin{matrix} 3 \\ \text{s}(m) \text{?s}(m) \\ (1) \end{matrix} \end{aligned}$$

defines a polynomial in

Empirical observations on the inverse moment matrix SOS polynomial

The inverse moment-matrix SOS polynomial is associated to a measure μ which satisfies the following. Assumption 1 μ is a Borel probability measure on \mathbb{R}^p with all its moments finite and $M_d(\mu)$ is positive definite for a given $d \in \mathbb{N}$. Definition 1 Let μ, d satisfy Assumption 1. We call the SOS polynomial $Q_{\mu,d}^{(2d)}$ defined by the application: $x \mapsto$

$$Q_{\mu,d}^{(2d)}(x) := v_d(x)^T M_d(\mu)^{-1} v_d(x),$$

$$x \in \mathbb{R}^p,$$

(2)

For the null polynomial, we use the convention that its degree is 0 and it is smaller than all other monomials.

3

the inverse moment-matrix SOS polynomial of degree $2d$ associated to μ . Actually, connection to orthogonal polynomials will show that the inverse function $x \mapsto Q_{\mu,d}^{(2d)}(x)$ is called the Christoffel function in the literature [16, 5] (see also Section 4). In the remainder of this section, we focus on the situation when μ corresponds to an empirical measure μ_P over n points in \mathbb{R}^p which are fixed. So let $x_1, \dots, x_n \in \mathbb{R}^p$ be a fixed set of points and let $\mu_P := \sum_{i=1}^n \frac{1}{n} \delta_{x_i}$ where δ_x corresponds to the Dirac measure at x . In such a case the polynomial $Q_{\mu,d}^{(2d)}$ in (2) is determined only by the empirical moments up to degree $2d$ of our collection of points. Note that we also require that $M_d(\mu) \succ 0$. In other words, the points x_1, \dots, x_n do not belong to an algebraic set defined by a polynomial of degree less or equal to d . We first describe empirical properties of inverse moment matrix SOS polynomial in this context of empirical measures. A mathematical intuition and further properties behind these observations are developed in Section 4. 3.1

Sublevel sets

The starting point of our investigations is the following phenomenon which to the best of our knowledge has remained unnoticed in the literature. For the sake of clarity and simplicity we provide an illustration in the plane. Consider the following experiment in \mathbb{R}^2 for a fixed $d \in \mathbb{N}$: represent on the same graphic, the cloud of points $\{x_i\}_{i=1}^n$ and the sublevel sets of SOS polynomial $Q_{\mu,d}^{(2d)}$ in \mathbb{R}^2 (equivalently, the superlevel sets of the Christoffel function). This is illustrated in the left panel of Figure 3. The collection of points consists of 500 simulations of two different Gaussians and the value of d is 4. The striking feature of this plot is that the level sets capture the global shape of the cloud of points quite accurately. In particular, the level set $\{x : Q_{\mu,d}^{(2d)}(x) \geq p+d \cdot d\}$ captures most of the points. We could reproduce very similar observations on different shapes with various number of points in \mathbb{R}^2 and degree d (see Appendix A). 3.2

Measuring outlyingness

An additional remark in a similar line is that $Q_{\mu,d}^{(2d)}$ tends to take higher values on points which are isolated from other points. Indeed in the left panel of Figure 3, the value of the polynomial tends to be smaller on the boundary of the cloud. This extends to situations where the collection of points correspond to shape with a high density of points with a few additional outliers. We reproduce a similar experiment on the right panel of Figure 3. In this example, 1000

points are sampled close to a ring shape and 40 additional points are sampled uniformly on a larger square. We do not represent the sublevel sets of $Q_{\gamma,d}$ here. Instead, the color and shape of the points are taken proportionally to the value of $Q_{\gamma,d}$, with $d = 8$. First, the results confirm the observation of the previous paragraph, points that fall close to the ring shape tend to be smaller and points on the boundary of the ring shape are larger. Second, there is a clear increase in the size of the points that are relatively far away from the ring shape. This highlights the fact that $Q_{\gamma,d}$ tends to take higher value in less populated areas of the space.

3.3 Relation to maximum likelihood estimation

If we fix $d = 1$, we recover the maximum P_n likelihood estimation P_n for the Gaussian, up to a constant additive factor. To see this, set $\gamma = n^{-1} \sum_{i=1}^n x_i x_i^T$ and $S = n^{-1} \sum_{i=1}^n x_i x_i^T$. With this notation, we have the following block representation of the moment matrix,

$$\begin{pmatrix} 1 + \gamma^T V \gamma & \gamma^T V \\ \gamma^T V^T & V \end{pmatrix} = \begin{pmatrix} M_d(\gamma) & \gamma^T S \\ \gamma S^T & V \end{pmatrix} = \begin{pmatrix} M_d(\gamma) & \gamma^T S \\ \gamma S^T & V \end{pmatrix}$$

where $V = S^{-1} \gamma^T \gamma$ is the empirical covariance matrix and the expression for the inverse is given by Schur complement. In this case, we have $Q_{\gamma,1}(x) = 1 + (x - \gamma)^T V^{-1} (x - \gamma)$ for all $x \in \mathbb{R}^p$. We recognize the quadratic form that appears in the density function of the multivariate Gaussian with parameters estimated by maximum likelihood. This suggests a connection between the inverse SOS moment polynomial and maximum likelihood estimation. Unfortunately, this connection is difficult to generalize for higher values of d and we do not pursue the idea of interpreting the empirical observations of this section through the prism of maximum likelihood estimation and leave it for further research. Instead, we propose an alternative view in Section 4.

3.4 Computational aspects

Recall that $s(d) = p+d$ is the number of p -variate monomials of degree up to d . The computation of $Q_{\gamma,d}$ requires $O(n s(d)^2)$ operations for the computation of the moment matrix and $O(s(d)^3)$ operations for the matrix inversion. The evaluation of $Q_{\gamma,d}$ requires $O(s(d)^2)$ operations. Estimating the coefficients of $Q_{\gamma,d}$ has a computational cost that depends only linearly in the number of points n . The cost of evaluating $Q_{\gamma,d}$ is constant with respect to the number of points n . This is an important contrast with kernel based or distance based methods (such as nearest neighbors and one class SVM) for density estimation or outlier detection since they usually require at least $O(n^2)$ operations for the evaluation of the model [3]. Moreover, this is well suited for online settings where inverse moment matrix computation can be done using rank one Woodbury updates [15, Section 2.7.1]. The dependence in the dimension p is of the order of p^d for a fixed d . Similarly, the dependence in d is of the order of d^p for a fixed dimension p and the joint dependence is exponential. Furthermore, $M_d(\gamma)$ has a Hankel structure which is known to produce ill conditioned matrices. This suggests that the direct computation and evaluation of $Q_{\gamma,d}$ will mostly make sense for moderate dimensions and degree d . In our experiments, for large d , the evaluation of $Q_{\gamma,d}$ remains quite stable, but the inversion leads to numerical error for higher values (around 20).

Invariance and interpretation through orthogonal polynomials

The purpose of this section is to provide a mathematical rationale that explains the empirical observations made in Section 3. All the proofs are postponed to Appendix B. We fix a Borel probability measure μ on \mathbb{R}^p which satisfies Assumption 1. Note that $M_d(\mu)$ is always positive definite if μ is not supported on the zero set of a polynomial of degree at most d . Under Assumption 1, $M_d(\mu)$ induces an inner product on $\mathcal{R}_s(d)$ and by extension on $\mathcal{R}_d[X]$ (see Section 2). This inner product is denoted by $\langle \cdot, \cdot \rangle_d$ and satisfies for any polynomials $P, Q \in \mathcal{R}_d[X]$ with coefficients $p, q \in \mathcal{R}_s(d)$, $\langle \sum p_i X^i, \sum q_i X^i \rangle_d = \sum p_i q_i = \langle p, q \rangle_d$.

We will also use the canonical inner product over $\mathcal{R}_d[X]$ which we write $\langle P, Q \rangle_d := \langle p, q \rangle_d$ for any polynomials $P, Q \in \mathcal{R}_d[X]$ with coefficients $p, q \in \mathcal{R}_s(d)$. We will omit the subscripts for this canonical inner product and use $\langle \cdot, \cdot \rangle$ for both products. 4.1

Affine invariance

It is worth noticing that the mapping $x \mapsto Q_{\mu,d}(x)$ does not depend on the particular choice of $\text{vd}(X)$ as a basis of $\mathcal{R}_d[X]$, any other basis would lead to the same mapping. This leads to the result that $Q_{\mu,d}$ captures affine invariant properties of μ . Lemma 1 Let μ satisfy Assumption 1 and $A \in \mathbb{R}^{p \times p}$, $b \in \mathbb{R}^p$ define an invertible affine mapping on \mathbb{R}^p , $A : x \mapsto Ax + b$. Then, the push forward measure, defined by $\mu'(S) = \mu(A^{-1}(S))$ for all Borel sets $S \subset \mathbb{R}^p$, satisfies Assumption 1 (with the same d as μ) and for all $x \in \mathbb{R}^p$, $Q_{\mu',d}(x) = Q_{\mu,d}(Ax + b)$. Proof Lemma 1 is probably better understood when $\mu = 1/n \sum_{i=1}^n \delta_{x_i}$ as in Section 3. In this case, we have $\mu' = 1/n \sum_{i=1}^n \delta_{Ax_i + b}$ and Lemma 1 asserts that the level sets of $Q_{\mu',d}$ are simply the images of those of $Q_{\mu,d}$ under the affine transformation $x \mapsto Ax + b$. This is illustrated in Appendix D. 4.2

Connection with orthogonal polynomials

We define a classical [16, 5] family of orthonormal polynomials, $\{P_i\}_{i=0}^{\infty}$ ordered according to \deg which satisfies for all $i \geq 0$ $\langle P_i, P_j \rangle = \delta_{ij}$, $\langle P_i, X \rangle = 0$ if $i \geq 1$, $\langle P_0, X \rangle = 0$ if $i \geq 1$, $\langle P_i, X \rangle = 0$ if $i \geq 1$.

(3)

It follows from (3) that $\langle P_i, P_j \rangle = 0$ if $i \neq j$. Existence and uniqueness of such a family is guaranteed by the Gram-Schmidt orthonormalization process following the \deg order, and by the

positivity of the moment matrix, see for instance [5, Theorem 3.1.11]. There exist determinantal formulae [9] and more precise description can be made for measures which have additional geometric properties, see [5] for many examples. Let $D_d(\mu)$ be the lower triangular matrix whose rows are the coefficients of the polynomials P_i defined in (3) ordered by \deg . It can be shown that $D_d(\mu) = L_d(\mu)^T$, where $L_d(\mu)$ is the Cholesky factorization of $M_d(\mu)$. Furthermore, there is a direct relation with the inverse moment matrix as $M_d(\mu)^{-1} = D_d(\mu)^T D_d(\mu)$ [9, Proof of Theorem 3.1]. This has the following consequence. Lemma 2 Let μ satisfy Assumption 1, then $Q_{\mu,d} = P_d^T$, where the family $\{P_i\}_{i=0}^{\infty}$ is \deg ordered and \mathcal{R}_p $Q_{\mu,d}(x) = s(d)$. That is, $Q_{\mu,d}$

is a very specific and distinguished SOS polynomial, the sum of squares of the orthonormal basis elements $\{P_i\}_{i=1}^N$ of $\mathcal{R}_d(X)$ (w.r.t. μ). Furthermore, the average value of Q_d with respect to μ is $s(d)$ which corresponds to the red level set in left panel of Figure 3.

A variational formulation for the inverse moment matrix SOS polynomial

In this section, we show that the family of polynomials $\{P_i\}_{i=1}^N$ defined in (3) is the unique d solution (up to a multiplicative constant) of a convex optimization problem over polynomials. This fact combined with Lemma 2 provides a mathematical rationale for the empirical observations outlined in Section 3. Consider the following optimization problem. $\min_{P \in \mathcal{R}_d(X)} \int P^2(x) d\mu(x)$ (4) $Q_i \succeq 0, \sum_{i=1}^N Q_i = P^2$ s.t. $Q_i \succeq 0, \int Q_i d\mu = 0, \int P d\mu = 0$

where $Q_i(x) = \sum_{j=1}^N q_{ij}(x)$ is a polynomial and μ is a real variable for each $i = 1, \dots, N$. We first comment on problem (4). Let $P = \sum_{i=1}^N Q_i$ be the SOS polynomial appearing in the objective function of (4). The objective of (4) simply involves the average value of P with respect to μ . Let $\mathcal{S}_d \subset \mathcal{R}_d[X]$ be the set of such SOS polynomials P which have a sum of square decomposition satisfying the constraints of (4) (for some arbitrary value of the real variables $\{q_{ij}\}_{i,j=1}^N$). With this \mathcal{R} notation, problem (4) has the simple formulation $\min_{P \in \mathcal{S}_d} \int P^2 d\mu$. Based on this formulation, problem (4) can be interpreted as balancing two antagonist targets. On one hand the minimization of the average value of the SOS polynomial P with respect to μ , on the other hand the avoidance of the trivial polynomial, enforced by the constraint that $P \in \mathcal{S}_d$. The constraint $P \in \mathcal{S}_d$ is simple and natural. It ensures that P is a sum of squares of polynomials $\{Q_i\}_{i=1}^N$, where d the leading term of each Q_i (according to the ordering μ) is $q_{ii}(x)$ with $q_{ii} \not\equiv 0$ (and hence does not vanish). Inversely, using Cholesky factorization, for any SOS polynomial Q of degree $2d$ which coefficient matrix (see equation (1)) is positive definite, there exists a $\lambda > 0$ such that $\lambda Q \in \mathcal{S}_d$. This suggests that \mathcal{S}_d is a quite general class of nonvanishing SOS polynomials. The following result, which gives a relation between Q_d and solutions of (4), uses a generalization of [16, Theorem 3.1.2] to several orthogonal polynomials of several variables. Theorem 1 : Under Assumption 1, problem (4) is a convex optimization problem with a unique optimal solution (Q_{opt}, μ_{opt}) , which satisfies $Q_{opt} \in \mathcal{S}_d$, for some $\mu_{opt} \not\equiv 0$. In particular, $\mu_{opt} \in \mathcal{P}_d$ the distinguished SOS polynomial $Q_d = P^2 = \sum_{i=1}^N (Q_{opt})^2$, is (part of) the unique d optimal solution of (4). Theorem 1 states that up to the scaling factor μ_{opt} , the distinguished SOS polynomial Q_d is the unique optimal solution of problem (4). A detailed proof is provided in the Appendix B and we only sketch the main ideas here. First, it is remarkable that for each fixed $d \in \mathbb{N}$ (and again up to μ_{opt} is the unique optimal solution of the problem: $\min_{P \in \mathcal{R}_d(X)} \int P^2 d\mu$ to a scaling factor) the polynomial $P^2 = \sum_{i=1}^N (Q_{opt})^2$ is (part of) the unique d optimal solution of (4). This fact is well-known in the univariate case [16, Theorem 3.1.2] and does not seem to have been exploited in the literature, at

least for purposes similar to ours. So intuitively, P^2 should be as close to 0 as possible on the support of μ . Problem (4) has similar properties and the \mathcal{R}

constraint on the vector of weights γ enforces that, at an optimal solution, the contribution of $(Q_{\gamma,d}^2 - d)^2$ to the overall sum in the criterion is the same for all γ . Using Lemma 2 yields (up to a multiplicative constant) the polynomial $Q_{\gamma,d}$. Other constraints on γ would yield different weighted sum of the squares $P_{\gamma,d}^2$. This will be a subject of further investigations. To sum up, Theorem 1 provides a rationale for our observations. Indeed when solving (4), intuitively, $Q_{\gamma,d}$ should be close to 0 on average while remaining in a class of nonvanishing SOS polynomials. 4.4

Christoffel function and outlier detection

The following result from [5, Theorem 3.5.6] draws a direct connection between $Q_{\gamma,d}$ and the Christoffel function (the right hand side of (5)). Theorem 2 ([5]) Let Assumption 1 hold and let $z \in \mathbb{R}^p$ be fixed, arbitrary. Then Z

$$Z = \min_{x \in \mathbb{R}^p} \frac{Q_{\gamma,d}^2(z)}{Q_{\gamma,d}^2(x)} = \frac{P_{\gamma,d}(z)}{P_{\gamma,d}(x)} \quad (5)$$

\mathbb{R}^p

Theorem 2 provides a mathematical rationale for the use of $Q_{\gamma,d}$ for outlier or novelty detection purposes. Indeed, from Lemma 2 and equation (3), we have $Q_{\gamma,d}^2 \geq 1$ on \mathbb{R}^p . Furthermore, the solution of the minimization problem in (5) satisfies $P_{\gamma,d}(z)^2 = 1$ and $\gamma \cdot x \in \mathbb{R}^p : P_{\gamma,d}(x)^2 \leq 1$ (inequality). Hence, for high values of $Q_{\gamma,d}(z)$, the sublevel set $\{x \in \mathbb{R}^p : Q_{\gamma,d}(x) \leq Q_{\gamma,d}(z)\}$ (by Markov's

$x \in \mathbb{R}^p : P_{\gamma,d}(x)^2 \leq 1$ contains most of the mass of γ while $P_{\gamma,d}(z)^2 = 1$. An illustration of this discussion is given in appendix E. Again the result of Theorem 2 does not seem to have been interpreted for purposes similar to ours.

5

Experiments on network intrusion datasets

In addition to having its own mathematical interest, Theorem 1 can be exploited for various purposes.

For instance, the sub-level sets of $Q_{\gamma,d}$, and in particular $\{x \in \mathbb{R}^p : Q_{\gamma,d}(x) \leq d\}$, can be used to encode a cloud of points in a simple and compact form. However in this section we focus on another potential application in anomaly detection. Empirical findings described in Section 3 suggest that the polynomial $Q_{\gamma,d}$ can be used to detect outliers in a collection of real vectors (with γ the empirical average). This is backed up by the results presented in Section 4. We illustrate these properties on a real world example. We choose the KDD cup 99 network intrusion dataset [13] consisting of network connection data, labeled as normal traffic or network intrusions. We follow [19] and [18] and construct five datasets consisting of labeled vectors in \mathbb{R}^3 with the following properties

http	567498	0.004
smtp	95156	0.0003
ftp-data	30464	0.023
ftp	4091	0.077
others	5858	0.016

The details on the datasets construction are available in [19, 18] and reproduced in Appendix C. The main idea is to compute an outlyingness score

(independant of the label) and compare outliers predicted by the score and network intrusion labels. The underlying assumption is that network intrusions correspond to infrequent abnormal behaviors and could be considered as outliers. We reproduce the same experiment as in [18, Section 5.4] using the value of Q^*,d from Definition 1 as an outlyingness score (with $d = 3$). The authors of [18] have compared different methods in the same experimental setting: robust estimation and Mahalanobis distance [8, 10], mixture models [14] and recurrent neural networks. The results are gathered in [18, Figure 7]. In the left panel of Figure 2 we represent the same performance measure for our approach: we first compute the value of Q^*,d for each datapoint and use it as an outlyingness score. We then display the proportion of correctly identified outliers, with score above a given threshold, as a function of the proportion of examples with score above the threshold (for different values of the threshold). The main comments are as follows. 7

```

d (AUPR)
dataset
0.8
http 0.6
smtp ftp_data
0.4
ftp
2 (0.18) 3 (0.18)
0.50
4 (0.16) 5 (0.15)
0.25
others
0.2
1 (0.08)
0.75
Precision
% correctly identified outliers
1.00
1.0
6 (0.13)
0.0
0.00 0.0
0.2
0.4
0.6
0.8
1.0
0.0
% top outlyingness score
0.2
0.4
0.6

```

0.8
1.0
Recall

Figure 2: Left: reproduction of the results described in [18] with the evaluation of $Q_{\gamma,d}$ as an outlyingness score ($d = 3$). Right: precision-recall curves for different values of d (dataset `others`). The inverse moment matrix SOS polynomial does detect network intrusions with varying performances on the five datasets. Except for the `ftp-data` dataset, the global shape of these curves are very similar to results reported in [18, Figure 7] indicating that the proposed approach is comparable to other dedicated methods for intrusion detection in these four datasets. In a second experiment, we investigate the effect of changing the value of d on the performances. We focus on the `others` dataset because it is the most heterogeneous. We adopt a slightly different measure of performance and use precision recall (see for example [4]) to measure performances in identifying network intrusions (the higher the curve, the better). We call the area under such curves the AUPR. The right panel of Figure 2 represents these results. First, the case $d = 1$, which corresponds to vanilla Mahalanobis distance as outlined in Section 3.3, gives poor performances. Second, the global performances rapidly increase with d and then decrease and stabilize. This suggests that d can be used as a tuning parameter to control the ‘complexity’ of $Q_{\gamma,d}$. Indeed, $2d$ is the degree of the polynomial $Q_{\gamma,d}$ and it is expected that more complex models will identify more diverse classes of examples as outliers. In our case, this means identifying regular traffic as outliers while it actually does not correspond to intrusions. In general, a good heuristic regarding the tuning of d is to investigate performances on a well specified task in a preliminary experiment.

6

Future work

An important question is the asymptotic regime when $d \rightarrow \infty$. Current state of knowledge suggests that, up to a correct scaling, the limit of the Christoffel functions (when known to exist) involves an edge effect term, related to the support of the measure, and the density of γ with respect to Lebesgue measure, see for example [2] for the Euclidean ball. It also suggests connections with the notion of equilibrium measure in potential theory [17, 1, 7]. Generalization and interpretation of these results in our context will be investigated in future work. Even though good approximations are obtained with low degree (at least in dimension 2 or 3), the approach involves the inversion of large ill conditioned Hankel matrices which reduces considerably the applicability for higher degrees and dimensions. A promising research line is to develop approximation procedures and advanced optimization and algebra tools so that the approach could scale computationally to higher dimensions and degrees. Finally, we did not touch the question of statistical accuracy. In the context of empirical processes, this will be very relevant to understand further potential applications in machine learning and reduce the gap between the abstract orthogonal polynomial theory and practical machine learning applications. Acknowledgments This work was partly supported by project ERC-ADG TAMING 666981, ERC-Advanced

Grant of the European Research Council and grant number FA9550-15-1-0500 from the Air Force Office of Scientific Research, Air Force Material Command.

8

2 References

- [1] R. J. Berman (2009). Bergman kernels for weighted polynomials and weighted equilibrium measures of C_n . *Indiana University Mathematics Journal*, 58(4):1921-1946.
- [2] L. Bos, B. Della Vecchia and G. Mastroianni (1998). On the asymptotics of Christoffel functions for centrally symmetric weights functions on the ball in R^n . *Rendiconti del Circolo Matematico di Palermo*, 52:277-290.
- [3] V. Chandola, A. Banerjee and V. Kumar (2009). Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41(3):15.
- [4] J. Davis and M. Goadrich (2006). The relationship between Precision-Recall and ROC curves. *Proceedings of the 23rd international conference on Machine learning* (pp. 233-240). ACM.
- [5] C.F. Dunkl and Y. Xu (2001). *Orthogonal polynomials of several variables*. Cambridge University Press. MR1827871.
- [6] G.H Golub, P. Milanfar and J. Varah (1999). A stable numerical method for inverting shape from moments. *SIAM Journal on Scientific Computing* 21(4):1222-1243 (1999).
- [7] B. Gustafsson, M. Putinar, E. Saff and N. Stylianopoulos (2009). Bergman polynomials on an archipelago: estimates, zeros and shape reconstruction. *Advances in Mathematics* 222(4):1405-1460.
- [8] A.S. Hadi (1994). A modification of a method for the detection of outliers in multivariate samples. *Journal of the Royal Statistical Society. Series B (Methodological)*, 56(2):393-396.
- [9] 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.
- [10] E.M. Knorr, R.T. Ng and R.H. Zamar (2001). Robust space transformations for distance-based operations. *Proceedings of the international conference on Knowledge discovery and data mining* (pp. 126-135). ACM.
- [11] J.B. Lasserre (2015). Level Sets and NonGaussian Integrals of Positively Homogeneous Functions. *International Game Theory Review*, 17(01):1540001.
- [12] J.B. Lasserre and M. Putinar (2015). Algebraic-exponential Data Recovery from Moments. *Discrete & Computational Geometry*, 54(4):993-1012.
- [13] M. Lichman (2013). UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml> University of California, Irvine, School of Information and Computer Sciences.
- [14] J.J. Oliver, R.A. Baxter and C.S. Wallace (1996). Unsupervised learning using MML. *Proceedings of the International Conference on Machine Learning* (pp. 364-372).
- [15] W. H. Press, S. A. Teukolsky, W. T. Vetterling and B. P. Flannery (2007). *Numerical Recipes: The Art of Scientific Computing* (3rd Edition). Cambridge University Press.
- [16] G. Szegő (1974). *Orthogonal polynomials*. In *Colloquium publications*, AMS, (23), fourth edition.
- [17] V. Totik (2000). Asymptotics for Christoffel functions for general measures on the real line. *Journal d'Analyse Mathématique*, 81(1):283-303.
- [18] G. Williams, R. Baxter, H. He, S. Hawkins and L. Gu (2002). A Comparative Study of RNN for Outlier Detection in Data Mining. *IEEE International Conference on Data Mining* (p. 709). IEEE

Computer Society. [19] K. Yamanishi, J.I. Takeuchi, G. Williams and P. Milne (2004). On-line unsupervised outlier detection using finite mixtures with discounting learning algorithms. *Data Mining and Knowledge Discovery*, 8(3):275-300.

9