The Christoffel-Darboux kernel for Data Analysis

Jean B. Lasserre

LAAS-CNRS, Institute of Mathematics & ANITI, Toulouse, France

Journées NUsCAP, Toulouse, 2022

Research funded by ANITI institute, under grant ANR-19-PIA3-0004 and ANR-NUsCAP-20-CE48-0014



Jean B. Lasserre semidefinite characterization

3

・ 同 ト ・ ヨ ト ・ ヨ ト …

Joint work with

- Edouard Pauwels (IRIT & IMT, Toulouse)
- Mihai Putinar (UCSB, USA & University of Newcastle UK)

個 とく ヨ とく ヨ とう

ambridge Monographs on Applied and Computational Mathematics

The Christoffel-Darboux Kernel for Data Analysis

Jean Bernard Lasserre, Edouard Pauwels and Mihai Putinar



Jean B. Lasserre semidefinite characterization

ъ



Jean B. Lasserre semidefinite characterization

Motivation

Consider the following cloud of 2D-points (data set) below



The red curve is the level set

$$egin{array}{lll} {old S}_{oldsymbol{\gamma}} \, := \, \{ \, {old X} : \, \, {old Q}_{oldsymbol{d}}({old X}) \leq \, {old \gamma} \, \}, \quad {old \gamma} \in \mathbb{R}_+ \, . \end{array}$$

of a certain polynomial $Q_d \in \mathbb{R}[x_1, x_2]$ of degree 2*d*.

Provide that S_{γ} captures quite well the shape of the cloud.

Not a coincidence!

Surprisingly, low degree *d* for Q_d is often enough to get a pretty good idea of the shape of Ω (at least in dimension p = 2, 3)



Jean B. Lasserre

semidefinite characterization

Cook up your own convincing example

Perform the following simple operations on a preferred cloud of 2*D*-points: So let d = 2, p = 2 and $s(d) = \binom{p+d}{p}$.

- Let $\mathbf{v}_d(\mathbf{x})^T = (1, x_1, x_2, x_1^2, x_1 x_2, \dots, x_1 x_2^{d-1}, x_2^d)$. be the vector of all monomials $x_1^i x_2^j$ of total degree $i + j \le d$
- Form the real symmetric matrix of size *s*(*d*)

$$\mathbf{M}_d := \frac{1}{N} \sum_{i=1}^N \mathbf{v}_d(\mathbf{x}(i)) \, \mathbf{v}_d(\mathbf{x}(i))^T \, ,$$

where the sum is over all points $(\mathbf{x}(i))_{i=1...,N} \subset \mathbb{R}^2$ of the data set.

(個) (日) (日) (日)

So the matrix $N \cdot \mathbf{M}_d$ reads:

$$\sum_{i} \begin{bmatrix} 1 & x_{1}(i) & x_{2}(i) & x_{1}(i)^{2} & \dots & x_{2}(i)^{d} \\ x_{1}(i) & x_{1}(i)^{2} & x_{1}(i) x_{2}(i) & x_{1}(i)^{3} & \dots & x_{1}(i) x_{2}(i)^{d} \\ x_{2}(i) & x_{1}(i) x_{2}(i) & x_{2}(i)^{2} & x_{1}(i)^{2} x_{2}(i) & \dots & x_{2}(i)^{d+1} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ x_{2}(i)^{d} & x_{1}(i) x_{2}(i)^{d} & x_{2}(i)^{d+1} & x_{1}(i)^{2} x_{2}(i)^{d} & \dots & x_{2}(i)^{2d} \end{bmatrix}$$

Jean B. Lasserre semidefinite characterization

◆□ > ◆□ > ◆臣 > ◆臣 > ─臣 ─のへで

 \mathbb{P} Note that typically, \mathbf{M}_d is what is called the MOMENT-matrix of the empirical measure

$$u^{\mathsf{N}} := \frac{1}{N} \sum_{i=1}^{N} \delta_{\mathbf{x}(i)}$$

associated with a sample of size N, drawn according to an unknown measure μ .

For the (usual) notation $\delta_{\mathbf{x}(i)}$ stands for the DIRAC measure supported at the point $\mathbf{x}(i)$ of \mathbb{R}^2 .

Next, form the SOS polynomial:

$$\mathbf{x} \mapsto Q_d(\mathbf{x}) := \mathbf{v}_d(\mathbf{x})^T \mathbf{M}_d^{-1} \mathbf{v}_d(\mathbf{x}).$$
$$= (\mathbf{1}, x_1, x_2, x_1^2, \dots, x_2^d) \mathbf{M}_d^{-1} \begin{pmatrix} \mathbf{1} \\ x_1 \\ x_2 \\ x_1^2 \\ \dots \\ x_2^d \end{pmatrix}$$

Plot some level sets

$$old S_\gamma := \{\, {f x} \in {\mathbb R}^2: \; {old Q}_d({f x}) \, = \, \gamma \, \}$$

for some values of γ , the thick one representing the particular value $\gamma = \binom{2+d}{2}$.

< 回 > < 回 > < 回 > … 回

The Christoffel function $\Lambda_d : \mathbb{R}^p \to \mathbb{R}_+$ is the reciprocal

 $\mathbf{x} \mapsto \mathbf{Q}_d(\mathbf{x})^{-1}, \quad \forall \mathbf{x} \in \mathbb{R}^p$

of the SOS polynomial Q_d .

It has a rich history in Approximation theory and Orthogonal Polynomials.

Among main contributors: Nevai, Totik, Króo, Lubinsky, Simon, ...

Image: Image:

・ 同 ト ・ ヨ ト ・ ヨ ト …

The Christoffel function $\Lambda_d : \mathbb{R}^{\rho} \to \mathbb{R}_+$ is the reciprocal

 $\mathbf{x} \mapsto \mathbf{Q}_{\mathbf{d}}(\mathbf{x})^{-1}, \quad \forall \mathbf{x} \in \mathbb{R}^{p}$

of the SOS polynomial Q_d .

It has a rich history in Approximation theory and Orthogonal Polynomials.

Among main contributors: Nevai, Totik, Króo, Lubinsky, Simon, ...

IF ... The CF seems to be not so well-known in data analysis

くぼう くほう くほう

Let μ be a (positive) measure supported on a compact set $\Omega \subset \mathbb{R}^p$ with nonempty interior.

A family
$$(\mathcal{P}_{\alpha})_{\alpha \in \mathbb{N}^{p}} \subset \mathbb{R}[\mathbf{x}]$$
 is orthonormal w.r.t. μ if
$$\int_{\Omega} \mathcal{P}_{\alpha}(\mathbf{x}) \, \mathcal{P}_{\beta}(\mathbf{x}) \, \mu(d\mathbf{x}) = \, \delta_{\alpha = \beta} \,, \quad \forall \alpha, \beta \in \mathbb{N}^{p} \,.$$

For Here $\delta_{\alpha=\beta}$ is the standard Kronecker symbol

> < 三 > < 三 > <</p>

How to construct a family $(P_{\alpha})_{\alpha \in \mathbb{N}^{p}}$

Let
$$\mathbb{N}_t^{p} := \{ \alpha \in \mathbb{N}^{p} : \sum_i \alpha_i \leq t \}$$
 and suppose that all moments

$$\boldsymbol{\mu}_{\alpha} := \int_{\Omega} \mathbf{x}^{\alpha} \, \boldsymbol{d} \boldsymbol{\mu} \,, \quad \forall \alpha \in \mathbb{N}_{2t}^{\boldsymbol{p}} \,,$$

are available.

For the one may construct an orthonormal family $(P_{\alpha})_{\alpha \in \mathbb{N}_{t}^{\rho}}$ from determinants of moment matrices associated with μ .

The moment matrix $\mathbf{M}_{d}(\mu)$ is the real symmetric matrix with rows and columns indexed by $(\mathbf{x}^{\alpha})_{\alpha \in \mathbb{N}_{d}^{p}}$, and with entries

$$\mathbf{M}_{d}(\mu)(\alpha,\beta) := \int_{\Omega} \mathbf{x}^{\alpha+\beta} \, d\mu = \mu_{\alpha+\beta}, \quad \forall \alpha,\beta \in \mathbb{N}_{d}^{p}.$$

Illustrative example in dimension 2:

$$\mathbf{M}_{1}(\mu) := \begin{pmatrix} 1 & X_{1} & X_{2} \\ 1 & \mu_{00} & \mu_{10} & \mu_{01} \\ X_{1} & \mu_{10} & \mu_{20} & \mu_{11} \\ X_{2} & \mu_{01} & \mu_{11} & \mu_{02} \end{pmatrix}$$

is the moment matrix of μ of "degree d=1".

通 とく ヨ とく ヨ とう

One way to construct polynomials orthonormal w.r.t. μ

Fix an ordering of \mathbb{N}^{p} (e.g. lexicographic ordering)

$$\underbrace{(0,0)}_{\text{degree0}}, \underbrace{(1,0), (0,1)}_{\text{degree1}}, \underbrace{(2,0), (1,1), (0,2)}_{\text{degree2}}, (3,0), (2,1), \dots$$
Then $P_{00}(\mathbf{x}) = 1$ for all $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$.
$$Q_{10}(\mathbf{x}) := \det \begin{pmatrix} \mu_{00} & \mu_{10} \\ 1 & X_1 \end{pmatrix} = X_1 - \mu_{10}.$$

$$Q_{01}(\mathbf{x}) := \det \begin{pmatrix} \mu_{00} & \mu_{10} \\ \mu_{10} & \mu_{20} & \mu_{11} \\ 1 & X_1 & X_2 \end{pmatrix}$$

 $= \mu_{10}\mu_{11} - \mu_{01}\mu_{20} - X_1 \left(\mu_{00}\mu_{11} - \mu_{10}\mu_{01}\right) + X_2 \left(\mu_{00}\mu_{20} - \mu_{10}^2\right)$

For Then normalize, i.e. $P_{10} = \theta Q_{10}$ with θ such that

$$heta^2\int_\Omega Q_{10}^2\, d\mu\,=\,1\,.$$

and similarly with $P_{01} = \theta Q_{01}$

One way to construct polynomials orthonormal w.r.t. μ

Fix an ordering of \mathbb{N}^{p} (e.g. lexicographic ordering)

$$\underbrace{(0,0)}_{\text{degree0}},\underbrace{(1,0),(0,1)}_{\text{degree1}},\underbrace{(2,0),(1,1),(0,2)}_{\text{degree2}},(3,0),(2,1),\dots$$
Then $P_{00}(\mathbf{x}) = 1$ for all $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$.
$$Q_{10}(\mathbf{x}) := \det \begin{pmatrix} \mu_{00} & \mu_{10} \\ 1 & X_1 \end{pmatrix} = X_1 - \mu_{10}.$$

$$Q_{01}(\mathbf{x}) := \det \begin{pmatrix} \mu_{00} & \mu_{10} \\ \mu_{10} & \mu_{20} & \mu_{11} \\ 1 & X_1 & X_2 \end{pmatrix}$$

 $= \mu_{10}\mu_{11} - \mu_{01}\mu_{20} - X_1 \left(\mu_{00}\mu_{11} - \mu_{10}\mu_{01}\right) + X_2 \left(\mu_{00}\mu_{20} - \mu_{10}^2\right)$

Then normalize, i.e. $P_{10} = \theta Q_{10}$ with θ such that

$$heta^2 \int_{\Omega} Q_{10}^2 \, d\mu \, = \, 1 \, .$$

and similarly with $P_{01} = \theta Q_{01}$.

Similarly,

$$Q_{20}(\mathbf{x}) := \det \begin{pmatrix} \mu_{00} & \mu_{10} & \mu_{01} & \mu_{20} \\ \mu_{10} & \mu_{20} & \mu_{11} & \mu_{30} \\ \mu_{01} & \mu_{11} & \mu_{02} & \mu_{21} \\ 1 & X_1 & X_2 & X_1^2 \end{pmatrix}$$

$$= X_1^2 \det \begin{pmatrix} \mu_{00} & \mu_{10} & \mu_{01} \\ \mu_{10} & \mu_{20} & \mu_{11} \\ \mu_{01} & \mu_{11} & \mu_{02} \end{pmatrix} - X_2 (\cdots) + X_1 (\cdots) - (\cdots).$$

and $P_{20} = \theta Q_{20}$ with θ such that

$$heta^2\int_{\Omega}Q_{20}^2\,d\mu\,=\,1\,.$$

◆□ > ◆□ > ◆臣 > ◆臣 > ─臣 ─のへで

The vector space $\mathbb{R}[\mathbf{x}]_d$ viewed as a subspace of $L^2(\mu)$ is a Reproducing Kernel Hilbert Space (RKHS). Its *reproducing kernel*

$$(\mathbf{x},\mathbf{y})\mapsto {\it K}^{\mu}_{\it d}(\mathbf{x},\mathbf{y})\,:=\,\sum_{|lpha|\leq {\it d}}{\it P}_{lpha}(\mathbf{x})\,{\it P}_{lpha}(\mathbf{y})\,,\quad orall\,\mathbf{x},\mathbf{y}\in\mathbb{R}^{\it p}\,,$$

is called the *Christoffel-Darboux kernel*.

프 > 프

The reproducing property

$$\mathbf{x}\mapsto q(\mathbf{x})\,=\,\int_{\Omega} {\mathcal K}^{\mu}_{d}(\mathbf{x},\mathbf{y})\,q(\mathbf{y})\,d\mu(\mathbf{y})\,,\quad orall q\in \mathbb{R}[\mathbf{x}]_{d}\,.$$

 \mathbb{P} useful to determinate the best degree-*d* $L^2(\mu)$ -polynomial approximation

0

$$\inf_{q\in\mathbb{R}[\mathbf{x}]_d}\|f-q\|_{L^2(\boldsymbol{\mu})}$$

of $f \in L^2(\mu)$. Indeed:

$$\mathbf{x} \mapsto \widehat{f_d}(\mathbf{x}) := \sum_{\alpha \in \mathbb{N}_d^p} (\overbrace{\int_{\Omega} f(y) P_{\alpha}(y) d\mu}^{\widetilde{f_{d,\alpha}}}) P_{\alpha}(\mathbf{x}) \in \mathbb{R}[\mathbf{x}]_d$$
$$= \arg \min_{q \in \mathbb{R}[\mathbf{x}]_d} ||f - q||_{L^2(\mu)}$$

or, equivalently: $\begin{aligned} & \int_{\Omega} (f - \widehat{f_d})^2 \, d\mu \to 0 \quad \text{as } d \to \infty \\ & \lim_{d \to \infty} \|f - \widehat{f_d}\|_{L^2(\mu)} = 0. \end{aligned}$

◆□ > ◆□ > ◆臣 > ◆臣 > ─臣 ─のへで

Recall that the support Ω of μ is compact with nonempty interior, and let $(P_{\alpha})_{\alpha \in \mathbb{N}^{p}}$ be a family of orthonormal polynomials w.r.t. μ .

Theorem

The Christoffel function $\Lambda^{\mu}_{d} : \mathbb{R}^{p} \to \mathbb{R}_{+}$ is defined by:

$$\xi\mapsto \Lambda^\mu_d(\xi)^{-1}\,=\,\sum_{|lpha|\leq d} P_lpha(\xi)^2\,=\, K^\mu_d(\xi,\xi)\,,\quad orall\xi\in\mathbb{R}^p\,,$$

and it also satisfies the variational property:

$$\Lambda^{\mu}_{d}(\xi) = \min_{P \in \mathbb{R}[\mathbf{x}]_{d}} \left\{ \int_{\Omega} P^{2} d\mu : P(\xi) = 1 \right\}, \quad \forall \xi \in \mathbb{R}^{p}.$$

Alternatively

$$\Lambda^{\mu}_{\boldsymbol{d}}(\xi)^{-1} \,=\, \mathbf{v}_{\boldsymbol{d}}(\xi)^T \mathbf{M}_{\boldsymbol{d}}(\mu)^{-1} \, \mathbf{v}_{\boldsymbol{d}}(\xi) \,, \quad \forall \xi \in \mathbb{R}^{\rho} \,.$$

くぼう くほう くほう

Recall that the support Ω of μ is compact with nonempty interior, and let $(P_{\alpha})_{\alpha \in \mathbb{N}^{p}}$ be a family of orthonormal polynomials w.r.t. μ .

Theorem

The Christoffel function $\Lambda^{\mu}_{d} : \mathbb{R}^{p} \to \mathbb{R}_{+}$ is defined by:

$$\xi\mapsto \Lambda^\mu_d(\xi)^{-1}\,=\,\sum_{|lpha|\leq d} P_lpha(\xi)^2\,=\, K^\mu_d(\xi,\xi)\,,\quad orall\xi\in\mathbb{R}^p\,,$$

and it also satisfies the variational property:

$$\Lambda^{\mu}_{d}(\xi) = \min_{P \in \mathbb{R}[\mathbf{x}]_{d}} \left\{ \int_{\Omega} P^{2} d\mu : P(\xi) = 1 \right\}, \quad \forall \xi \in \mathbb{R}^{p}.$$

Alternatively

$$\Lambda^{\mu}_{d}(\xi)^{-1} \,=\, \mathbf{v}_{d}(\xi)^{T} \mathbf{M}_{d}(\mu)^{-1} \, \mathbf{v}_{d}(\xi) \,, \quad \forall \xi \in \mathbb{R}^{p}$$

The support of the support Ω of the underlying measure μ .

Theorem

Let the support Ω of μ be compact with nonempty interior. Then:

- For all $\mathbf{x} \in \operatorname{int}(\Omega)$: $K_d^{\mu}(\mathbf{x}, \mathbf{x}) = O(d^p)$.
- For all $\mathbf{x} \in int(\mathbb{R}^p \setminus \Omega)$: $K_d^{\mu}(\mathbf{x}, \mathbf{x}) = \Omega(\exp(\alpha d))$ for some $\alpha > 0$.

```
In particular, as d \to \infty,
d^{\rho} \Lambda^{\mu}_{d}(\mathbf{x}) \to 0 very fast whenever \mathbf{x} \notin \Omega.
```

・ 同 ト ・ ヨ ト ・ ヨ ト

Growth rates for $K_d^{\mu}(\mathbf{x}, \mathbf{x}) = \Lambda_d^{\mu}(\mathbf{x})^{-1}$.



◆□> ◆□> ◆豆> ◆豆> ・豆 ・ のへで

• Under some (restrictive) assumption on Ω and μ

$$\lim_{d\to\infty} s(d) \Lambda^{\mu}_{d}(\xi) = f_{\mu}(\xi) \omega(\xi)^{-1}$$

where ω is the density of an equilibrium measure intrinsically associated with Ω . For instance with p = 1 and $\Omega = [-1, 1]$, $\omega(\xi) = \sqrt{1 - \xi^2}$.

If μ and ν have same support Ω and respective densities f_μ and f_ν w.r.t. Lebesgue measure on Ω, positive on Ω, then:

$$\lim_{d\to\infty}\frac{\Lambda^{\mu}_{d}(\xi)}{\Lambda^{\nu}_{d}(\xi)} = \frac{f_{\mu}(\xi)}{f_{\nu}(\xi)}, \quad \forall \xi \in \Omega.$$

useful for density approximation

通 とくほ とくほ とう

• Under some (restrictive) assumption on Ω and μ

$$\lim_{d\to\infty} s(d) \Lambda^{\mu}_{d}(\xi) = f_{\mu}(\xi) \omega(\xi)^{-1}$$

where ω is the density of an equilibrium measure intrinsically associated with Ω . For instance with p = 1 and $\Omega = [-1, 1]$, $\omega(\xi) = \sqrt{1 - \xi^2}$.

If μ and ν have same support Ω and respective densities f_μ and f_ν w.r.t. Lebesgue measure on Ω, positive on Ω, then:

$$\lim_{d\to\infty}\frac{\Lambda^{\mu}_{d}(\xi)}{\Lambda^{\nu}_{d}(\xi)} = \frac{f_{\mu}(\xi)}{f_{\nu}(\xi)}, \quad \forall \xi \in \Omega.$$

useful for density approximation

If Ω is not full-dimensional and is supported on a real variety $V \subset \mathbb{R}^{p}$, then for sufficiently large degree *d*:

 $d \mapsto \operatorname{rank}(\mathbf{M}_d) = q(d)$

where $q \in \mathbb{R}[t]$ is the Hilbert polynomial associated with *V* and whose degree provides the dimension of *V*.

So one may use the rank of the moment matrix \mathbf{M}_d to identify the dimension of the underlying variety.

we useful for manifold learning.

🗇 🕨 🖉 🖻 🖉 🖉 🖉

The Christoffel function can be used in several important applications of Machine Learning (e.g. outlier detection, density approximation, manifold learning). In this case the measure μ is the empirical probability measure μ^N associated with a cloud of N points $\mathcal{C} \subset \mathbb{R}^p$ (the data of interest).

The computing $\Lambda_d^{\mu^N}$ requires only one pass over the data & no optimization

御下 《唐下 《唐下 》 唐

The Christoffel function can be used in several important applications of Machine Learning (e.g. outlier detection, density approximation, manifold learning). In this case the measure μ is the empirical probability measure μ^N associated with a cloud of N points $\mathcal{C} \subset \mathbb{R}^p$ (the data of interest).

Computing $\Lambda_d^{\mu^N}$ requires only one pass over the data & no optimization

< 回 > < 回 > < 回 > … 回

The Christoffel function can be used in several important applications of Machine Learning (e.g. outlier detection, density approximation, manifold learning). In this case the measure μ is the empirical probability measure μ^N associated with a cloud of N points $\mathcal{C} \subset \mathbb{R}^p$ (the data of interest).

Computing $\Lambda_d^{\mu^N}$ requires only one pass over the data & no optimization

< 回 > < 回 > < 回 > … 回

Rank-one update

Updating the Christoffel function when the cloud of N points one additional point ξ is added to the cloud of N points is easy.

$$(N+1)\mu^{N+1} = \sum_{i=1}^{N} \delta_{\mathbf{x}(i)} + \delta_{\boldsymbol{\xi}} = N\mu^{N} + \delta_{\boldsymbol{\xi}}$$

By Sherman-Morrison's rank-one update formula

$$((N+1) \mathbf{M}_{d}(\mu^{N+1}))^{-1} = (N \mathbf{M}_{d}(\mu^{N}) + \mathbf{v}_{d}(\xi) \mathbf{v}_{d}(\xi)^{T})^{-1} = (N \mathbf{M}_{d}(\mu^{N}))^{-1} - \frac{1}{N^{2}} \frac{\mathbf{M}_{d}(\mu^{N})^{-1} \mathbf{v}_{d}(\xi) \mathbf{v}_{d}(\xi)^{T} \mathbf{M}_{d}(\mu^{N})^{-1}}{1 + \mathbf{v}_{d}(\xi) \mathbf{M}_{d}(\mu^{N})^{-1} \mathbf{v}_{d}(\xi)}$$

ъ

and therefore

one obtains the simple update formula:

$$\frac{1}{N+1} \Lambda_{d}^{\mu^{N+1}}(\mathbf{x}) = \frac{1}{N} \left[\Lambda_{d}^{\mu^{N}}(\mathbf{x}) - \frac{K_{d}^{\mu^{N}}(x,\xi)^{2}}{N(1+\Lambda_{d}^{\mu^{N}}(\mathbf{x}))} \right], \quad \forall \mathbf{x}$$
$$\frac{1}{N+1} \Lambda_{d}^{\mu^{N+1}}(\xi) = \frac{1}{N} \Lambda_{d}^{\mu^{N}}(\xi) - \frac{1}{N^{2}} \frac{\Lambda_{d}^{\mu^{N}}(\xi)^{2}}{1+\Lambda_{d}^{\mu^{N}}(\xi)}$$

▲□▶ ▲□▶ ▲目▶ ▲目▶ 三目 のへで

For instance one may decide to classify as outliers all points $\boldsymbol{\xi}$ such that $\Lambda_d^{\mu N}(\boldsymbol{\xi}) < {p+d \choose p}^{-1}$.

Such a strategy (even with relatively low degree d) is as efficient as more elaborated techniques, with only one parameter (the degree d), and with no optimization involved.

Lass. & Pauwels (2016) Sorting out typicality via the inverse moment matrix SOS polynomial, NIPS 2016.
 Lass. & Pauwels (2019) The empirical Christoffel function with applications in data analysis, Adv. Comp. Math. 45, pp. 1439–1468

・ 同 ト ・ ヨ ト ・ ヨ ト

For instance one may decide to classify as outliers all points $\boldsymbol{\xi}$ such that $\Lambda_d^{\mu N}(\boldsymbol{\xi}) < {p+d \choose p}^{-1}$.

Such a strategy (even with relatively low degree *d*) is as efficient as more elaborated techniques, with only one parameter (the degree *d*), and with no optimization involved.

Lass. & Pauwels (2016) Sorting out typicality via the inverse moment matrix SOS polynomial, NIPS 2016.
 Lass. & Pauwels (2019) The empirical Christoffel function with applications in data analysis, Adv. Comp. Math. 45, pp. 1439–1468

(雪) (ヨ) (ヨ)

For instance one may decide to classify as outliers all points $\boldsymbol{\xi}$ such that $\Lambda_d^{\mu N}(\boldsymbol{\xi}) < {p+d \choose p}^{-1}$.

Such a strategy (even with relatively low degree *d*) is as efficient as more elaborated techniques, with only one parameter (the degree *d*), and with no optimization involved.

Lass. & Pauwels (2016) Sorting out typicality via the inverse moment matrix SOS polynomial, NIPS 2016.
 Lass. & Pauwels (2019) The empirical Christoffel function with applications in data analysis, Adv. Comp. Math. 45, pp. 1439–1468

(雪) (ヨ) (ヨ)

Promising results in a recent collaboration with:

1- S. Dauzère-Péres, V. Borodin (EMSE) and the STMicroelectronics company for

data analysis of processing times for operations in a job-shop (e.g. detection of anomalies, density estimation, etc.)

2- L. Travé, K. Ducharlet (LAAS-CNRS) and the Carl Berger-Levrault company for detection of anomalies in data analysis of wireless sensors network used in several applications (e.g. units of air treatment, automatic bagage conveyor in airports (data in form of temporal series), IFF

K. Ducharlet, L. Travé, J.B. Lasserre, M.V. Le Lann, Y. Miloudi. Leveraging the Christoffel Function for Outlier Detection in Data Streams, submitted.

◎ ▶ ★ 臣 ▶ ★ 臣 ▶

Promising results in a recent collaboration with:

1- S. Dauzère-Péres, V. Borodin (EMSE) and the STMicroelectronics company for

data analysis of processing times for operations in a job-shop (e.g. detection of anomalies, density estimation, etc.)

2- L. Travé, K. Ducharlet (LAAS-CNRS) and the Carl Berger-Levrault company for detection of anomalies in data analysis of wireless sensors network used in several applications (e.g. units of air treatment, automatic bagage conveyor in airports (data in form of temporal series),

K. Ducharlet, L. Travé, J.B. Lasserre, M.V. Le Lann, Y. Miloudi. Leveraging the Christoffel Function for Outlier Detection in Data Streams, submitted. A measure μ on compact set Ω is completely determined by its moments and therefore it should not be a surprise that its moment matrix $\mathbf{M}_d(\mu)$ contains a lot of information.

We have already seen that its inverse $M_d(\mu)^{-1}$ defines the Christoffel function.

When μ is degenerate and its support Ω is contained in a real algebraic variety then the kernel of $\mathbf{M}_d(\mu)$ identifies the generators of a corresponding ideal of $\mathbb{R}[\mathbf{x}]$.

(雪) (ヨ) (ヨ)

For instance let $\Omega \subset \mathbb{S}^{p-1}$ (the Euclidean unit sphere of \mathbb{R}^p)



3

ヘロン 人間 とくほ とくほ とう

Then the kernel of $\mathbf{M}_d(\mu)$ contains vectors of coefficients of polynomials in the ideal generated by the quadratic polynomial $\mathbf{x} \mapsto g(\mathbf{x}) := 1 - \|\mathbf{x}\|^2$.

In fact and remarkably,

 $\operatorname{rank} \mathbf{M}_d(\mu) = p(d)$

for some univariate polynomial p (the Hilbert polynomial associated with the algebraic variety) which is of degree t if t is the dimension of the variety.

For instance t = p - 1 if the support is contained in the sphere \mathbb{S}^{p-1} of \mathbb{R}^{p} .

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

For $\varepsilon > 0$ sufficiently small, the Christoffel function

$$\mathbf{x} \mapsto \Lambda^{\boldsymbol{\mu}}_{d}(\mathbf{x}) \,=\, \mathbf{v}_{d}(\mathbf{x}) \, (\mathbf{M}_{d}(\boldsymbol{\mu}) + \varepsilon \, \boldsymbol{I})^{-1} \, \mathbf{v}_{d}(\mathbf{x})$$

and its empirical version (from a sample of data points on Ω)

$$\mathbf{x} \mapsto \Lambda_d^{\mu^{N}}(\mathbf{x}) = \mathbf{v}_d(\mathbf{x}) \left(\mathbf{M}_d(\mu^{N}) + \varepsilon \mathbf{I} \right)^{-1} \mathbf{v}_d(\mathbf{x})$$

identifies correctly the support of Ω .

Pauwels E., Putinar M., Lass. J.B. (2021). Data analysis from empirical moments and the Christoffel function, Found. Comput. Math. 21, pp. 243–273. Again this illustrates how quite sophisticated concepts of algebraic geometry are hidden and encapsulated in the moment matrix $M_d(\mu)$.

They can be exploited to extract various useful information on the data set.

In addition, extraction of this information can be done via quite simple linear algebra techniques. Again this illustrates how quite sophisticated concepts of algebraic geometry are hidden and encapsulated in the moment matrix $M_d(\mu)$.

They can be exploited to extract various useful information on the data set.

In addition, extraction of this information can be done via quite simple linear algebra techniques. Again this illustrates how quite sophisticated concepts of algebraic geometry are hidden and encapsulated in the moment matrix $M_d(\mu)$.

They can be exploited to extract various useful information on the data set.

In addition, extraction of this information can be done via quite simple linear algebra techniques.

However

for non modest dimension of data, matrix inversion of \mathbf{M}_d^{-1} does not scale well ...

¹²⁷ On the other hand

for evaluation $\Lambda^{\mu}_{d}(\xi)$ at a point $\xi \in \mathbb{R}^{p}$, the variational formulation

$$\Lambda^{\mu}_{d}(\xi) = \min_{P \in \mathbb{R}[\mathbf{x}]_{d}} \left\{ \int_{\Omega} P^{2} d\mu : P(\xi) = 1 \right\}, \quad \forall \xi \in \mathbb{R}^{p}.$$

is the simple quadratic programming problem.

$$\min_{\boldsymbol{p}\in\mathbb{R}^{s(d)}} \{ \boldsymbol{p}^T \mathbf{M}_d \boldsymbol{p} : \mathbf{v}_d(\boldsymbol{\xi})^T \boldsymbol{p} = 1 \},\$$

which can be solved quite efficiently.

However

for non modest dimension of data, matrix inversion of \mathbf{M}_d^{-1} does not scale well ...

On the other hand

for evaluation $\Lambda^{\mu}_{d}(\xi)$ at a point $\xi \in \mathbb{R}^{p}$, the variational formulation

$$\Lambda^{\mu}_{d}(\boldsymbol{\xi}) = \min_{\boldsymbol{P} \in \mathbb{R}[\mathbf{x}]_{d}} \left\{ \int_{\Omega} \boldsymbol{P}^{2} \, d\mu : \, \boldsymbol{P}(\boldsymbol{\xi}) = 1 \right\}, \quad \forall \boldsymbol{\xi} \in \mathbb{R}^{p}.$$

is the simple quadratic programming problem.

$$\min_{\boldsymbol{p}\in\mathbb{R}^{s(d)}} \{ \boldsymbol{p}^T \mathbf{M}_d \boldsymbol{p} : \mathbf{v}_d(\boldsymbol{\xi})^T \boldsymbol{p} = \mathbf{1} \},\$$

which can be solved quite efficiently.

< 🗇 > < 🖻 > .

Other non-polynomial kernels, some popular in ML (e.g. Gaussian kernels), can be very efficient, to provide a large class of functions on which efficient calculation in large dimension is possible. However they are not related (at least directly) to an underlying measure supported on the data points.

Real Again, a distinguishing feature of the CD-kernel is its deep connexion with the underlying measure.

- It not only "encodes" the cloud of data points,
- but it also captures many essential features of the more complex measure supported on those data points.

Should be seen as another item in the arsenal of kernel methods in ML.

(画) (ヨ) (ヨ)

Other non-polynomial kernels, some popular in ML (e.g. Gaussian kernels), can be very efficient, to provide a large class of functions on which efficient calculation in large dimension is possible. However they are not related (at least directly) to an underlying measure supported on the data points.

- Again, a distinguishing feature of the CD-kernel is its deep connexion with the underlying measure.
 - It not only "encodes" the cloud of data points,
 - but it also captures many essential features of the more complex measure supported on those data points.

Should be seen as another item in the arsenal of kernel methods in ML.

(雪) (ヨ) (ヨ)

Other non-polynomial kernels, some popular in ML (e.g. Gaussian kernels), can be very efficient, to provide a large class of functions on which efficient calculation in large dimension is possible. However they are not related (at least directly) to an underlying measure supported on the data points.

- Again, a distinguishing feature of the CD-kernel is its deep connexion with the underlying measure.
 - It not only "encodes" the cloud of data points,
 - but it also captures many essential features of the more complex measure supported on those data points.

Should be seen as another item in the arsenal of kernel methods in ML.

伺き くほき くほう

II: The CF to approximate piecewise continuous functions.

A typical approach is to approximate $f : [0, 1] \rightarrow \mathbb{R}$ in some function space, e.g. its projection on $\mathbb{R}[\mathbf{x}]_n \subset L^2([0, 1])$:

$$x\mapsto \hat{f}_n(x) := \sum_{j=0}^n \left(\int_0^1 f(y) L_j(y) dy\right) L_j(x),$$

with an orthonormal basis $(L_j)_{j \in \mathbb{N}}$ of $L^2([0, 1])$.



Ex: Chebyshev interpolant

Typical Gibbs phenomenon occurs.

Jean B. Lasserre semidefinite characterization

Alternative Positive Kernels with better convergence properties have been proposed, still in the same framework:

Féjer, Jackson kernels, etc.

- Reproducing property of the CD kernel is LOST
- Preserve positivity (e.g when approximating a density)
- Better convergence properties than the CD kernel, in particular uniform convergence (for continuous functions) on arbitrary compact subsets

Observe that

$$\hat{f}_n = \int_0^1 K_n^\lambda(x, y) f(y) \, dy \, ,$$

where K_n^{λ} is the CD-kernel of Lebesgue λ on [0, 1].

A counter-intuitive detour: Instead of considering $f : [0, 1] \rightarrow \mathbb{R}$

Solution $\Omega \subset \mathbb{R}^2$ of *f*, i.e., the set

 $\Omega := \{ (x, f(x)) : x \in [0, 1] \}.$

and the measure $d\phi(x, y) = \delta_{f(x)}(dy) dx$ supported on Ω .

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ …

Why should we do that as it implies going to \mathbb{R}^2 instead of staying in $\mathbb{R}?$

🖙 ... because

The support of φ is exactly the graph of f, and
The CF (x, y) → Λ^φ_n(x, y) identifies the support of φ!

▲□ ▶ ▲ ■ ▶ ▲ ■ ▶ ■ ● ● ● ●

$$\mathbf{v}_d(x, y) := (1, x, y, x^2, x y, y^2, \dots, x y^{d-1}, y^d).$$

Compute the degree-d empirical moment matrix:

$$\mathbf{M}_{d} := \sum_{i=1}^{N} \mathbf{v}_{d}((x_{i}, f(x_{i})) \mathbf{v}_{d}(x_{i}, f(x_{i}))^{T},$$

by one pass over the data

Compute the Christoffel function

$$x \mapsto \Lambda_d(x, y)^{-1} := \mathbf{v}_d(x, y)^T \mathbf{M}_d^{-1} \mathbf{v}_d(x, y)$$

Approximate f(x) by f_d(x) := arg min_y ∧_d(x, y)⁻¹.
 Image minimize a univariate polynomial! (easy)

$$\mathbf{v}_d(x, y) := (1, x, y, x^2, x y, y^2, \dots, x y^{d-1}, y^d).$$

Compute the degree-d empirical moment matrix:

$$\mathbf{M}_{d} := \sum_{i=1}^{N} \mathbf{v}_{d}((x_{i}, f(x_{i})) \mathbf{v}_{d}(x_{i}, f(x_{i}))^{T},$$

by one pass over the data

Compute the Christoffel function

$$x \mapsto \Lambda_d(x,y)^{-1} := \mathbf{v}_d(x,y)^T \mathbf{M}_d^{-1} \mathbf{v}_d(x,y).$$

Approximate f(x) by f_d(x) := arg min_y ∧_d(x, y)⁻¹.
 I[™] minimize a univariate polynomial! (easy)

$$\mathbf{v}_d(x, y) := (1, x, y, x^2, x y, y^2, \dots, x y^{d-1}, y^d).$$

Compute the degree-d empirical moment matrix:

$$\mathbf{M}_{d} := \sum_{i=1}^{N} \mathbf{v}_{d}((x_{i}, f(x_{i})) \mathbf{v}_{d}(x_{i}, f(x_{i}))^{T},$$

by one pass over the data

Compute the Christoffel function

$$x \mapsto \Lambda_d(x, y)^{-1} := \mathbf{v}_d(x, y)^T \mathbf{M}_d^{-1} \mathbf{v}_d(x, y).$$

• Approximate f(x) by $\hat{f}_d(x) := \arg \min_y \Lambda_d(x, y)^{-1}$.

$$\mathbf{v}_d(x,y) := (1, x, y, x^2, x y, y^2, \dots, x y^{d-1}, y^d).$$

Compute the degree-d empirical moment matrix:

$$\mathbf{M}_{d} := \sum_{i=1}^{N} \mathbf{v}_{d}((x_{i}, f(x_{i})) \mathbf{v}_{d}(x_{i}, f(x_{i}))^{T},$$

by one pass over the data

• Proceeding the Christoffel function

$$x \mapsto \Lambda_d(x,y)^{-1} := \mathbf{v}_d(x,y)^T \mathbf{M}_d^{-1} \mathbf{v}_d(x,y).$$

Approximate f(x) by f_d(x) := arg min_y ∧_d(x, y)⁻¹.
 № minimize a univariate polynomial! (easy)

Good convergence properties as $d \uparrow$

- \mathbb{E}^{1} -convergence,
- 🖙 even pointwise convergence on open sets with no point of discontinuity, and so almost uniform convergence.

S. Marx, E. Pauwels, T. Weisser, D. Henrion, J.B. Lass. Semi-algebraic approximation using Christoffel-Darboux kernel, Constructive Approximation, 2021

(個) (日) (日) 日



◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ● ●

Suppose that the function $f : [0, 1] \rightarrow \mathbb{R}$ to approximate is only known via its Fourier-Legendre coefficients

$$\phi_{i,1} = \int_0^1 x^i f(x) dx, \quad i = 0, 1, \dots$$

and we do not have access to other moments

$$\phi_{i,j} = \int_0^1 x^i f(x)^j dx, \quad j > 1; \ i = 0, 1, \dots$$

of the measure $\phi(d(\mathbf{x}, y)) = \delta_{f(\mathbf{x})}(dy) \lambda(d\mathbf{x})$

個人 くほん くほん 一日

Recall that $\lambda = (\lambda_i)_{i \in \mathbb{N}}$ is the moment-sequence of Lebesgue measure on [0, 1], and consider the semidefinite programs indexed by $n \in \mathbb{N}$:

$$\begin{aligned} \mathbf{P}_n : & \inf_{\boldsymbol{\psi}} \quad \{ \Theta_n(\boldsymbol{\psi}) : \, \mathbf{M}_n(\boldsymbol{\psi}) \succeq 0 \\ & \psi_{i,0} = \lambda_i, \quad i \in \mathbb{N} \\ & \psi_{i,1} = \phi_{i,1}, \quad i \in \mathbb{N} \, \} \,, \end{aligned}$$

where the inf is over all pseudo-moments $\psi = (\psi_{i,j})_{i,j \in \mathbb{N}^2_{2n}}$, and Θ_n is a certain linear functional.

個人 くほん くほん しほ

For every measure ν on $[0, 1] \times \mathbb{R}$, let $\nu = (\nu_{i,j})_{i,j \in \mathbb{N}}$ be the sequence of its moments.

Theorem

(i) For every $n \in \mathbb{N}$,

 $\Theta_n(\phi) \leq \Theta_n(\nu),$

for all measures ν on $[0, 1] \times \mathbb{R}$ whose moment-sequence ν is a feasible solution of \mathbf{P}_n .

(ii) Let ψ^n be an optimal solution of \mathbf{P}_n . Then

$$\lim_{n \to \infty} \psi_{i,j}^n = \phi_{i,j} = \int_{[0,1]} x^i f(x)^j dx, \quad \forall i, j = 0, 1, \dots$$

For Hence one may approximate accurately from finitely moments $\phi_{i,j}$ as described earlier.

D. Henrion & J.B. Lass. Graph recovery from incomplete moment information (2021), Constructive Approximation.

Christoffel function and Positive polynomials

Let $\Omega \subset \mathbb{R}^n$ be the basic semi-algebraic set (with nonempty interior)

$$\Omega := \{ \mathbf{x} \in \mathbb{R}^n : g_j(\mathbf{x}) \ge 0, \quad j = 1, \dots, m \}$$

with $g_j \in \mathbb{R}[\mathbf{x}]_{d_i}$ and let $s_j = \lceil \deg(g_j)/2 \rceil$. Let $g_0 = 1$ with $s_0 = 0$.

With t fixed, its associated quadratic module

$$Q_t(\Omega) := \{ \sum_{j=0}^m \sigma_j g_j : \sigma_j \in \Sigma[\mathbf{x}]_{t-s_j} \} \subset \mathbb{R}[\mathbf{x}]$$

is a convex cone with nonempty interior,

▲□ > ▲ □ > ▲ □ > …

$$Q_t(\Omega)^* := \{ \mathbf{y} \in \mathbb{R}^{s(t)} : \mathbf{M}_{t-s_j}(\mathbf{g}_j \mathbf{y}) \succeq 0, \quad j = 0, \dots, m \},$$

where $s(t) = \binom{n+t}{n}.$

Notice that if $\mathbf{M}_t(\mathbf{y})^{-1} \succ 0$ for all t

wł

one may define a family of polynomials $(P_{\alpha})_{\alpha \in \mathbb{N}^n} \subset \mathbb{R}[\mathbf{x}]$ orthonormal w.r.t. *y*, meaning that

$$L_{\mathbf{y}}(\mathbf{P}_{\alpha}\cdot\mathbf{P}_{\beta}) = \delta_{\alpha=\beta}, \quad \alpha, \beta \in \mathbb{N}^{n},$$

and exactly as for measures, the Christoffel function Λ_t^y

$$\mathbf{x} \mapsto \Lambda^{\mathbf{y}}_t(\mathbf{x})^{-1} := \sum_{|\alpha| \leq t} \mathcal{P}_{\alpha}(\mathbf{x})^2.$$

イロン 不良 とくほう 不良 とうしょう

Theorem

For every $p \in int(Q_t(\Omega))$ there exists $y \in int(Q_t(\Omega)^*)$ such that

$$p(\mathbf{x}) = \sum_{j=0}^{m} \left(\mathbf{v}_{t-s_j}(\mathbf{x})^T \mathbf{M}_t(g_j \mathbf{y})^{-1} \mathbf{v}_{t-s_j}(\mathbf{x}) \right) g_j(\mathbf{x})$$
$$= \sum_{j=0}^{m} \Lambda_{t-s_j}^{g_j \cdot \mathbf{y}}(\mathbf{x})^{-1} g_j(\mathbf{x})$$

where $(\boldsymbol{g} \cdot \boldsymbol{y})$ is the sequence of pseudo-moments

$$(\boldsymbol{g} \cdot \boldsymbol{y})_{\alpha} := \sum_{\gamma} \boldsymbol{g}_{\gamma} \, \boldsymbol{y}_{\alpha+\gamma}, \quad \alpha \in \mathbb{N}^{n} \quad (\text{if } \boldsymbol{g}(\mathbf{x}) = \sum_{\gamma} \boldsymbol{g}_{\gamma} \, \mathbf{x}^{\gamma}).$$

In addition $L_{\mathbf{y}}(\mathbf{p}) = \sum_{j=0}^{m} \binom{n+t-s_j}{n}$.

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ののの

The proof combines

- \mathbb{C} a result by Nesterov on a one-to-one correspondence between $int(Q_t(\Omega))$ and $int(Q_t(\Omega)^*)$, and

- 😰 the fact that

$$\mathbf{v}_{t-s_j}(\mathbf{x})^T \mathbf{M}_t(g_j \mathbf{y})^{-1} \mathbf{v}_{t-s_j}(\mathbf{x}) = \Lambda_{t-s_j}^{g_j \cdot \mathbf{y}}(\mathbf{x})^{-1}$$

|| (同) || (回) || (\cup) |

In other words:

In Putinar certificate

$$oldsymbol{
ho} \,=\, \sum_{j=0}^m \sigma_j \, g_j \,, \quad \sigma_j \in \mathbb{R}[\mathbf{x}]_{t-s_j} \,,$$

of strict positivity on Ω ,

 \mathbb{P} one may always choose the SOS weights σ_i in the form

$$\sigma_j(\mathbf{x}) := \Lambda_{t-\mathbf{s}_i}^{g_j \cdot \mathbf{y}}(\mathbf{x})^{-1}, \quad j = 0, \dots, m,$$

for some sequence of pseudo-moments $\mathbf{y} \in int(\mathbf{Q}_t(\Omega)^*)$.

(4回) (4回) (4回)

ъ

THANK YOU!

Jean B. Lasserre semidefinite characterization

◆□> ◆□> ◆豆> ◆豆> ・豆 ・ のへで