



HAL
open science

A hierarchy of convex relaxations for the total variation distance

Jean-Bernard Lasserre

► **To cite this version:**

Jean-Bernard Lasserre. A hierarchy of convex relaxations for the total variation distance. 2023. hal-04367575v1

HAL Id: hal-04367575

<https://laas.hal.science/hal-04367575v1>

Preprint submitted on 30 Dec 2023 (v1), last revised 20 Sep 2024 (v3)

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.

A HIERARCHY OF CONVEX RELAXATIONS FOR THE TOTAL VARIATION DISTANCE

JEAN B. LASSERRE

ABSTRACT. Given two measures μ, ν on \mathbb{R}^d that satisfy Carleman's condition, we provide a numerical scheme to approximate as closely as desired the total variation distance between μ and ν . It consists of solving a sequence (hierarchy) of convex relaxations whose associated sequence of optimal values converges to the total variation distance, an additional illustration of the versatility of the Moment-SOS hierarchy. Indeed each relaxation in the hierarchy is a semidefinite program whose size increases with the number of involved moments. It has an optimal solution which is a couple of degree- $2n$ pseudo-moments which converge, as n grows, to moments of the Hahn-Jordan decomposition of $\mu - \nu$.

1. INTRODUCTION

Evaluating a “distance” between measures has become an important topic with many applications, especially in Data Science and Machine Learning in particular. Among possible choices, the Wasserstein distance has become popular and one reason is that its optimal transport formulation allows to define efficient specialized procedures (e.g. the Sinkhorn algorithm) for its computation [10].

Contribution. In this paper we show that the *total variation* distance is also amenable to practical computation under fairly weak assumptions and so could provide an alternative to other distances when needed; see e.g. [9] for a discussion on relative merits of several distances. The total variation distance being the same as the Wasserstein distance with (nasty) cost function $c(x, y) = 1_{x \neq y}(x, y)$, is an indication that its effective computation is a computational challenge.

An important application of the total variation is in computer vision where it was introduced in [11] as a regularization criterion in some inverse problems (e.g. denoising of images), followed by [3]. Next, [2] inspired by [3], has proposed an algorithm to minimize the (discrete) total variation of functions $u \in L^1(\Omega)$ for a compact set $\Omega \subset \mathbb{R}^2$ (which in principle can be adapted to higher dimensions). (A discretized image is seen a $N \times N$ matrix.)

Our focus is different from the above cited works [2, 3, 11]. Independently of a particular application, and in a rather general context, we provide a numerical scheme to approximate as closely as desired the total variation distance between two measures μ and ν . We assume that all moments of μ and ν are finite, and that both

The author is supported by the AI Interdisciplinary Institute ANITI funding through the french program “Investing for the Future PISA” under the grant agreement number ANR-19-PI3A-0004. This research is also part of the programme DesCartes and is supported by the National Research Foundation, Prime Minister’s Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) programme.

μ and ν satisfy Carleman's condition. In particular we do *not* assume that μ or ν has compact support. We can formulate the problem as an instance of the Generalized Moment Problem and show that it is amenable to practical computation via the Moment-SOS hierarchy [7, 4]. As a result, one may approximate as closely as desired $\|\mu - \nu\|_{TV}$ as more and more moments of μ and ν are taken into account. Our contribution is to provide an additional tool in the arsenal of algorithms available in applied probability, for approximating as closely as desired, the total variation distance $\|\mu - \nu\|_{TV}$ based on moment information.

(i) This numerical scheme consists of solving a sequence (hierarchy) of convex relaxations. Each convex relaxation of the hierarchy is a semidefinite program¹ whose size increases with the number of moments of μ and ν involved.

(ii) The associated sequence of optimal values is monotone non decreasing and converges from below to $\|\mu - \nu\|_{TV}$.

(iii) At last but not least, the associated sequence of optimal solutions of relaxations (a couple of vectors whose size increases), converges to the couple of infinite moment-vectors of the Hahn-Jordan decomposition (ϕ_+^*, ϕ_-^*) of the signed measure $\mu - \nu$.

We wish to emphasize the weak assumption on the measures μ, ν , namely that they satisfy Carleman's condition (no compact support is required). It is a feature of the total variation distance to not discriminate between mutually singular probability measures (their distance is constant to 1). Relatively surprisingly, and as a good sign of the numerical scheme's behavior, it is shown in Example 1 that the exact distance $\|\delta_0 - \delta_\varepsilon\|_{TV}$ between the two (mutually singular) Dirac measures at 0 and $\varepsilon > 0$, is obtained at the first semidefinite relaxation of the hierarchy, irrespective of the value of $\varepsilon > 0$ (whereas one might have expected that convergence of the relaxations would depend on ε).

Interestingly, and as an alternative to algorithms based on a discretization (like e.g. Sinkhorn algorithm), the Wasserstein distance $W_2(\mu, \nu)$ (with *polynomial* cost $c(x, y)$) can also be approximated as closely as desired by a mesh-free practical computation by (i) applying the Moment-SOS hierarchy [4, 8] for solving the associated optimal transport problem (OT), and (ii) extract the transport map from the moment vector solution of the OT, by a non-standard application of the Christoffel-Darboux kernel [5]. However, crucial is the fact that the cost function is a polynomial (which of course excludes the nasty cost function $1_{x \neq y}(x, y)$).

2. MAIN RESULT

2.1. Notation and definitions. Let $\mathbb{R}[\mathbf{x}]$ denote the ring of real polynomials in the variables (x_1, \dots, x_d) and $\mathbb{R}[\mathbf{x}]_n \subset \mathbb{R}[\mathbf{x}]$ be its subset of polynomials of total degree at most n . Let $\mathbb{N}_n^d := \{\boldsymbol{\alpha} \in \mathbb{N}^d : \sum_i \alpha_i \leq n\}$ with cardinal $s(n) = \binom{n+d}{n}$. Let $\mathbf{v}_n(\mathbf{x}) = (x^\boldsymbol{\alpha})_{\boldsymbol{\alpha} \in \mathbb{N}_n^d}$ be the vector of monomials up to degree n , and let $\Sigma[\mathbf{x}]_n \subset \mathbb{R}[\mathbf{x}]_{2n}$ be the convex cone of polynomials of total degree at most $2n$ which are sum-of-squares (in short SOS). A polynomial $p \in \mathbb{R}[\mathbf{x}]_n$ can be identified with its vector of coefficients $\mathbf{p} = (p_\alpha) \in \mathbb{R}^{s(n)}$ in the monomial basis, and reads

$$\mathbf{x} \mapsto p(\mathbf{x}) := \langle \mathbf{p}, \mathbf{v}_n(\mathbf{x}) \rangle, \quad \forall p \in \mathbb{R}[\mathbf{x}].$$

¹A semidefinite program is a convex conic optimization problem that can be solved efficiently, up to arbitrary precision fixed in advance; see e.g. [1]

Denote by $\mathcal{M}(\mathbb{R}^d)$ (resp. $\mathcal{M}(\mathbb{R}^d)_+$) the space of signed (resp. positive) Borel measures on \mathbb{R}^d . For two Borel measures $\mu, \nu \in \mathcal{M}(\mathbb{R}^d)_+$, the notation $\mu \leq \nu$ stands for $\mu(B) \leq \nu(B)$ for all Borel sets $B \in \mathcal{B}(\mathbb{R}^d)$. The support of a Borel measure μ on \mathbb{R}^d is the smallest closed set A such that $\mu(\mathbb{R}^d \setminus A) = 0$, and such a set A is unique. A Borel measure whose all moments are finite is said to be (moment) *determinate* if there is no other measure with same moments.

For a real symmetric matrix $\mathbf{A} = \mathbf{A}^T$, the notation $\mathbf{A} \succeq 0$ (resp. $\mathbf{A} \succ 0$) stands for \mathbf{A} is positive semidefinite (p.s.d.) (resp. positive definite (p.d.)).

Hahn-Jordan decomposition. Given two finite Borel measures $\mu, \nu \in \mathcal{M}(\mathbb{R}^d)_+$, the signed measure $\mu - \nu$ has a unique Hahn-Jordan decomposition (ϕ_+^*, ϕ_-^*) such that $\phi_+^* - \phi_-^* = \mu - \nu$. That is, there exists a Borel set $A \in \mathcal{B}(\mathbb{R}^d)$ and two mutually singular positive measures ϕ_+^*, ϕ_-^* such that $\phi_+^*(\mathbb{R}^d) = \phi_+^*(A)$ while $\phi_-^*(A) = 0$, and

$$(2.1) \quad \phi_+^*(B) = (\mu - \nu)(B \cap A); \quad \phi_-^*(B) = (\nu - \mu)(B \cap (\mathbb{R}^d \setminus A)), \quad \forall B \in \mathcal{B}(\mathbb{R}^d).$$

In addition, and obviously, $\|\mu - \nu\|_{TV} \leq \mu(1) + \nu(1)$. Moreover, observe that $\phi_+^* \leq \mu$ and $\phi_-^* \leq \nu$.

Riesz linear functional and moment matrix. With a real sequence $\phi = (\phi_\alpha)_{\alpha \in \mathbb{N}^d}$ (in bold) is associated the *Riesz* linear functional $\phi \in \mathbb{R}[\mathbf{x}]^*$ (not in bold) defined by

$$p (= \sum_{\alpha} p_{\alpha} \mathbf{x}^{\alpha}) \mapsto \phi(p) = \langle \phi, \mathbf{p} \rangle = \sum_{\alpha} p_{\alpha} \phi_{\alpha}, \quad \forall p \in \mathbb{R}[\mathbf{x}],$$

and the moment matrix $\mathbf{M}_n(\phi)$ with rows and columns indexed by \mathbb{N}_n^d (hence of size $s(n)$), and with entries

$$\mathbf{M}_n(\phi)(\alpha, \beta) := \phi(\mathbf{x}^{\alpha+\beta}) = \phi_{\alpha+\beta}, \quad \alpha, \beta \in \mathbb{N}_n^d.$$

Notice that one may write indifferently $\mathbf{M}_n(\phi)$ or $\mathbf{M}_n(\phi)$, i.e., referring to the sequence ϕ truncated to degree- $2n$ moments or to the Riesz linear functional ϕ associated with ϕ .

A real sequence $\phi = (\phi_{\alpha})_{\alpha \in \mathbb{N}^d}$ has a representing measure if its associated linear functional ϕ is a Borel measure on \mathbb{R}^d . In this case $\mathbf{M}_n(\phi) \succeq 0$ for all n ; the converse is not true in general.

Carleman's condition. A sequence $\boldsymbol{\mu} = (\mu_{\alpha})_{\alpha \in \mathbb{N}^d}$ satisfies Carleman's condition if

$$(2.2) \quad \forall i = 1, \dots, d : \sum_{j=1}^{\infty} \mu(x_i^{2j})^{-1/2j} = +\infty.$$

The following theorem is due to Nussbaum:

Theorem 2.1. ([8, Theorem 3.5]) *Let a sequence $\boldsymbol{\mu} = (\mu_{\alpha})_{\alpha \in \mathbb{N}^d}$ be such that $\mathbf{M}_n(\boldsymbol{\mu}) \succeq 0$, for all $n \in \mathbb{N}$. If $\boldsymbol{\mu}$ satisfies Carleman's condition (2.2) then $\boldsymbol{\mu}$ has a representing measure μ on \mathbb{R}^d and μ is determinate.*

A sufficient condition to ensure that a measure μ satisfies the multivariate Carleman's condition is that

$$(2.3) \quad \int \exp(c|x_i|) d\mu < \infty, \quad i = 1, \dots, d,$$

if for some scalar $c > 0$.

2.2. A preliminary result.

Lemma 2.2. *Let $\mu, \varphi \in \mathcal{M}(\mathbb{R}^d)_+$ have finite moments and assume that μ satisfies Carleman's condition (2.2). Then*

$$(2.4) \quad \varphi \leq \mu \quad \Leftrightarrow \quad \mathbf{M}_n(\varphi) \preceq \mathbf{M}_n(\mu), \quad \forall n \in \mathbb{N}.$$

Proof. \Rightarrow is straightforward. Indeed:

$$\mu \geq \varphi \Rightarrow \left[\int p^2 d\mu \geq \int p^2 d\varphi, \forall p \in \mathbb{R}[\mathbf{x}] \right] \Rightarrow \mathbf{M}_n(\mu) \succeq \mathbf{M}_n(\varphi), \forall n \in \mathbb{N}.$$

\Leftarrow Assume that $\mathbf{M}_n(\varphi) \preceq \mathbf{M}_n(\mu)$ for all $n \in \mathbb{N}$, and consider the sequence $\gamma = (\gamma_\alpha)_{\alpha \in \mathbb{N}^d}$, with $\gamma_\alpha = \mu_\alpha - \varphi_\alpha$, for all $\alpha \in \mathbb{N}^d$. Then $\int x_i^{2n} d\varphi \leq \int x_i^{2n} d\mu$ for all n , and as Carleman's condition (2.2) holds for μ , we infer $\gamma(x_i^{2n}) \leq \mu(x_i^{2n})$ for all n , and all $i = 1, \dots, d$. This implies that γ satisfies Carleman's condition (2.2) and therefore, as $\mathbf{M}_n(\gamma) = \mathbf{M}_n(\mu) - \mathbf{M}_n(\varphi) \succeq 0$ for all n , we deduce that γ has a determinate representing measure γ on \mathbb{R}^d . In particular:

$$\int \mathbf{x}^\alpha d(\gamma + \varphi) = \gamma_\alpha + \varphi_\alpha = \mu_\alpha = \int \mathbf{x}^\alpha d\mu, \quad \forall \alpha \in \mathbb{N}^d \Rightarrow \gamma + \varphi = \mu,$$

where the last statement follows from determinateness of μ . Hence $\varphi \leq \mu$. \square

2.3. Main result. Given two finite Borel measures μ and ν on \mathbb{R}^d , introduce the infinite-dimensional LP:

$$(2.5) \quad \tau = \inf_{\phi^+, \phi^- \in \mathcal{M}(\mathbb{R}^d)_+} \{ \phi^+(1) + \phi^-(1) : \phi_+ - \phi_- = \mu - \nu \}.$$

Proposition 2.3. *The LP (2.5) has a unique optimal solution (ϕ_+^*, ϕ_-^*) which is the Hahn-Jordan decomposition of the signed measure $\mu - \nu$, and therefore $\tau = \phi_+^*(1) + \phi_-^*(1) = \|\mu - \nu\|_{TV}$.*

Proof. Let (ϕ^+, ϕ^-) be an arbitrary feasible solution of (2.5). Then as $\phi^+ - \phi^- = \mu - \nu$ one obtains $\phi^+(1) + \phi^-(1) \geq \|\phi^+ - \phi^-\|_{TV} = \|\mu - \nu\|_{TV}$. On the other hand, the Hahn-Jordan decomposition (ϕ_+^*, ϕ_-^*) of $\mu - \nu$ is feasible for (2.5), with value $\|\mu - \nu\|_{TV}$, whence the result. \square

Unfortunately the LP (2.5) is not very useful as it stands. It is just a particular rephrasing of the total variation distance between μ and ν . However we next see the a slight reinforcement of (2.5) will turn out to be very useful when passing to some hierarchy of convex relaxations. Indeed:

Proposition 2.4. *The linear program*

$$(2.6) \quad \inf_{\phi^+, \phi^- \in \mathcal{M}(\mathbb{R}^d)_+} \{ \phi^+(1) + \phi^-(1) : \phi_+ - \phi_- = \mu - \nu; \quad \phi^+ \leq \mu; \phi^- \leq \nu \}$$

has same optimal value $\tau = \|\mu - \nu\|_{TV}$, and optimal solution (ϕ_+^, ϕ_-^*) as (2.5).*

Proof. By construction, the optimal value ρ of (2.6) satisfies $\rho \geq \tau = \|\mu - \nu\|_{TV}$. On the other hand, with (ϕ_+^*, ϕ_-^*) being the Hahn-Jordan decomposition of $\mu - \nu$, observe that $\phi_+^* \leq \mu$, and $\phi_-^* \leq \nu$. Therefore (ϕ_+^*, ϕ_-^*) is an optimal solution of (2.6). Equivalently, the constraints $\phi^+ \leq \mu$ and $\phi^- \leq \nu$ are automatically satisfied at the optimal solution (ϕ_+^*, ϕ_-^*) of (2.5) and therefore (2.5) and (2.6) have same optimal value and same optimal solution. \square

Next, from now on we make the following assumption:

Assumption 2.5. (i) All moments of μ and ν are finite, and

(ii) μ and ν satisfy (2.3) (hence satisfy Carleman's condition (2.2)) for some scalar $c > 0$.

Consider the optimization problem

$$(2.7) \quad \hat{\tau} = \min_{\phi^+, \phi^- \in \mathcal{M}(\mathbb{R}^d)_+} \{ \phi^+(1) + \phi^-(1) : \phi^+ - \phi^- = \mu - \nu; \\ \mathbf{M}_n(\phi^+) \preceq \mathbf{M}_n(\mu); \mathbf{M}_n(\phi^-) \preceq \mathbf{M}_n(\nu), \quad \forall n \in \mathbb{N} \}.$$

Corollary 2.6. Let Assumption 2.5 hold. Then the Hahn-Jordan decomposition (ϕ_+^*, ϕ_-^*) of the signed measure $\mu - \nu$, is the unique optimal solution of (2.7), and $\hat{\tau} = \tau = \|\mu - \nu\|_{TV}$.

Proof. By Lemma 2.2, (2.6) and (2.7) are equivalent. \square

The nice feature of the LP (2.7) when compared to its equivalent formulation (2.6), is that the cost as well as the constraints of (2.7) can next be formulated in terms of moments of $(\mu, \nu, \phi^+, \phi^-)$, so as to yield the optimization problem:

$$(2.8) \quad \rho = \min_{\phi^+, \phi^- \in \mathcal{M}(\mathbb{R}^d)_+} \{ \phi^+(1) + \phi^-(1) : \\ \int \mathbf{x}^\alpha d(\phi^+ - \phi^-) = \int \mathbf{x}^\alpha d(\mu - \nu), \quad \forall \alpha \in \mathbb{N}^d \\ \mathbf{M}_n(\phi^+) \preceq \mathbf{M}_n(\mu); \mathbf{M}_n(\phi^-) \preceq \mathbf{M}_n(\nu), \quad \forall n \in \mathbb{N} \},$$

which is an instance of the Generalized Moment Problem (GMP); see e.g. [8].

Corollary 2.7. Let Assumption 2.5 hold. Then the Hahn-Jordan decomposition (ϕ_+^*, ϕ_-^*) of the signed measure $\mu - \nu$, is the unique optimal solution of (2.8), and $\rho = \|\mu - \nu\|_{TV}$.

Proof. Let (ϕ^+, ϕ^-) be an arbitrary feasible solution of (2.8). By Lemma 2.2, $\phi^+ \leq \mu$ and $\phi^- \leq \nu$. Hence $\phi^+ + \nu \leq \mu + \nu$, and $\phi^- + \mu \leq \mu + \nu$. As Assumption 2.5(ii) holds,

$$\int \exp(c|x_i|) d(\phi^+ + \nu) < \int \exp(c|x_i|) d(\mu + \nu) < \infty \\ \int \exp(c|x_i|) d(\phi^- + \mu) < \int \exp(c|x_i|) d(\mu + \nu) < \infty,$$

and therefore the measure $\phi^+ + \nu$ (resp. $\phi^- + \mu$) is determinate. But then the constraint $\int \mathbf{x}^\alpha d(\phi^+ - \phi^-) = \int \mathbf{x}^\alpha d(\mu - \nu)$ for all $\alpha \in \mathbb{N}^d$ reads:

$$\int \mathbf{x}^\alpha d(\phi^+ + \nu) = \int \mathbf{x}^\alpha d(\phi^- + \mu), \quad \forall \alpha \in \mathbb{N}^d,$$

which implies $\phi^+ + \nu = \phi^- + \mu$ by determinacy of the measures. Therefore (ϕ^+, ϕ^-) is a feasible solution of (2.7) with same value. In other words, (2.8) is equivalent to (2.7), whence the result. \square

2.4. A convergent hierarchy of semidefinite relaxations. As (2.8) is an instance of the GMP, it is natural to apply the Moment-SOS hierarchy [4, 7]. With each fixed $n \in \mathbb{N}$, consider the optimization problem

$$(2.9) \quad \rho_n = \min_{\phi, \psi} \left\{ \begin{array}{l} \phi(1) + \psi(1) : \quad \phi_{\alpha} - \psi_{\alpha} = \mu_{\alpha} - \nu_{\alpha}, \quad \forall \alpha \in \mathbb{N}_{2n}^d; \\ 0 \preceq \mathbf{M}_n(\phi) \preceq \mathbf{M}_n(\mu); \quad 0 \preceq \mathbf{M}_n(\psi) \preceq \mathbf{M}_n(\nu), \end{array} \right.$$

where now the optimization is over degree- $2n$ pseudo-moment vectors $\phi = (\phi_{\alpha})_{\alpha \in \mathbb{N}_{2n}^d}$ and $\psi = (\psi_{\alpha})_{\alpha \in \mathbb{N}_{2n}^d}$ (hence not necessarily coming from measures ϕ and ψ on \mathbb{R}^d). Of course (2.9) is an obvious relaxation of (2.8) and therefore $\rho_n \leq \rho = \|\mu - \nu\|_{TV}$ for all $n \in \mathbb{N}$.

Observe that for each fixed $n \in \mathbb{N}$, (2.9) is a semidefinite program that can be solved by off-the-shelf solvers like GloptiPoly [6] or Jump [12] (package of the Julia programming language).

Theorem 2.8. *Let Assumption 2.5 hold.*

(i) *For every fixed $n \in \mathbb{N}$, the optimization problem (2.9) has an optimal solution denoted $(\phi^{(n)}, \psi^{(n)})$.*

(ii) *In addition, $\rho_n \uparrow \|\mu - \nu\|_{TV}$ as $n \rightarrow \infty$, and moreover,*

$$(2.10) \quad \lim_{n \rightarrow \infty} \phi_{\alpha}^{(n)} = \int \mathbf{x}^{\alpha} d\phi_{+}^{*}; \quad \lim_{n \rightarrow \infty} \psi_{\alpha}^{(n)} = \int \mathbf{x}^{\alpha} d\phi_{-}^{*}, \quad \forall \alpha \in \mathbb{N}^d,$$

where $(\phi_{+}^{*}, \phi_{-}^{*})$ is the Hahn-Jordan decomposition of the signed measure $\mu - \nu$.

Proof. (i) Let (ϕ, ψ) be an arbitrary feasible solution of (2.9). As $\mathbf{M}_n(\phi) \preceq \mathbf{M}_n(\mu)$ one obtains

$$\phi(1) \leq \mu(1); \quad \phi(x_i^{2n}) \leq \mu(x_i^{2n}), \quad \forall i = 1, \dots, d,$$

and therefore, as $\mathbf{M}_n(\phi) \succeq 0$, by [8, Proposition 3.6],

$$(2.11) \quad |\phi_{\alpha}| \leq \max[\mu(1), \max_i \mu(x_i^{2d})], \quad \forall \alpha \in \mathbb{N}_{2n}^d.$$

Similarly, as $0 \preceq \mathbf{M}_n(\psi) \preceq \mathbf{M}_n(\nu)$,

$$(2.12) \quad |\psi_{\alpha}| \leq \max[\nu(1), \max_i \nu(x_i^{2d})], \quad \forall \alpha \in \mathbb{N}_{2n}^d.$$

Therefore the feasible set of (2.9) is compact. Hence (2.9) has an optimal solution.

(ii) For each fixed $n \in \mathbb{N}$, and since $\mathbf{M}_k(\phi^{(n)})$ is a submatrix of $\mathbf{M}_n(\phi^{(n)})$ for all $k = 1, \dots, n$, again by [8,],

$$\forall \alpha : 2k-1 \leq |\alpha| \leq 2k : |\phi_{\alpha}^{(n)}| \leq \max[\mu(1), \max_i \mu(x_i^{2k})] =: a_k; \quad k = 1, \dots, n,$$

and similarly

$$|\psi_{\alpha}^{(n)}| \leq \max[\nu(1), \max_i \nu(x_i^{2k})] =: b_k, \quad \forall \alpha : 2k-1 \leq |\alpha| \leq 2k; \quad k = 1, \dots, n.$$

Next, introduce the new infinite pseudo-moment sequences:

$$(2.13) \quad \hat{\phi}_{\alpha}^{(n)} := \phi_{\alpha}^{(n)} / a_k, \quad \forall \alpha : 2k-1 \leq |\alpha| \leq 2k; \quad k = 1, \dots, n,$$

and $\hat{\phi}_{\alpha}^{(n)} = 0$ for all $\alpha \in \mathbb{N}^d$ with $|\alpha| > 2n$. Similarly,

$$(2.14) \quad \hat{\psi}_{\alpha}^{(n)} := \psi_{\alpha}^{(n)} / b_k, \quad \forall \alpha : 2k-1 \leq |\alpha| \leq 2k; \quad k = 1, \dots, n,$$

and $\hat{\psi}_{\alpha}^{(n)} = 0$ for all $\alpha \in \mathbb{N}^d$ with $|\alpha| > 2n$.

Both sequences $\hat{\phi}^{(n)}$ and $\hat{\psi}^{(n)}$ are considered as elements of the unity ball $\mathbf{B}(0, 1)$ of the Banach space ℓ_∞ of uniformly bounded sequences, which is sequentially compact in the $\sigma(\ell_\infty, \ell_1)$ weak topology. Therefore there exist $\hat{\phi}, \hat{\psi} \in \mathbf{B}(0, 1)$ and a subsequence $(n_k)_{k \in \mathbb{N}}$ such that

$$(2.15) \quad \lim_{k \rightarrow \infty} \hat{\phi}_\alpha^{(n_k)} = \hat{\phi}_\alpha; \quad \lim_{k \rightarrow \infty} \hat{\psi}_\alpha^{(n_k)} = \hat{\psi}_\alpha, \quad \forall \alpha \in \mathbb{N}^d.$$

By doing the reverse scaling of (2.13)-(2.14), one obtains:

$$(2.16) \quad \forall \alpha \in \mathbb{N}^d : \quad \lim_{k \rightarrow \infty} \phi_\alpha^{(n_k)} = \phi_\alpha; \quad \lim_{k \rightarrow \infty} \psi_\alpha^{(n_k)} = \psi_\alpha,$$

where

$$\phi_\alpha := a_k \cdot \hat{\phi}_\alpha; \quad \psi_\alpha := b_k \cdot \hat{\psi}_\alpha; \quad \forall \alpha : 2k - 1 \leq |\alpha| \leq 2k; \quad k \in \mathbb{N}.$$

Fix $t \in \mathbb{N}$ arbitrary. As $\mathbf{M}_t(\phi^{(n)}) \succeq 0$ for all $n \geq t$, then by (2.16), $0 \preceq \mathbf{M}_t(\phi) \preceq \mathbf{M}_t(\mu)$, and as t was arbitrary, $0 \preceq \mathbf{M}_n(\phi) \preceq \mathbf{M}_n(\mu)$ for all n , and similarly $0 \preceq \mathbf{M}_n(\psi) \preceq \mathbf{M}_n(\nu)$ for all n .

Next, as $\mathbf{M}_n(\phi) \preceq \mathbf{M}_n(\mu)$, and μ satisfies Carleman's condition, then so does ϕ , and as $\mathbf{M}_n(\phi) \succeq 0$ for all n , it follows that $\phi = (\phi_\alpha)_{\alpha \in \mathbb{N}^d}$ has a representing measure ϕ on \mathbb{R}^d . Similarly, for same reasons, ψ has a representing measure ψ on \mathbb{R}^d .

In addition, by (2.16),

$$\|\mu - \nu\|_{TV} \geq \lim_{k \rightarrow \infty} \rho_{n_k} = \lim_{k \rightarrow \infty} \phi^{(n_k)}(1) + \psi^{(n_k)}(1) = \phi(1) + \psi(1),$$

and

$$\forall \alpha \in \mathbb{N}^d : \quad \mu_\alpha - \nu_\alpha = \lim_{k \rightarrow \infty} \phi_\alpha^{(n_k)} - \psi_\alpha^{(n_k)} = \phi_\alpha - \psi_\alpha.$$

Hence (ϕ, ψ) is an optimal solution of (2.7) (hence of (2.5) as well), and by Corollary 2.6, $(\phi, \psi) = (\phi_+^*, \phi_-^*)$, the Hahn-Jordan decomposition of $\mu - \nu$.

Finally, as the $(n_k)_{k \in \mathbb{N}}$ was an arbitrary converging subsequence and the limit is independent of the subsequence, the whole sequence converges. \square

If μ and ν are two probability measures, mutually singular, then $\|\mu - \nu\|_{TV} = 2$. A perfect case to check whether (2.9) is efficient, is to test (2.9) with the toy univariate example where $\mu = \delta_0$ and $\nu = \delta_\varepsilon$ for small value of $\varepsilon > 0$. Indeed, one might expect that the convergence $\rho_n \uparrow \|\mu - \nu\|_{TV}$ as n grows, could depend on ε (the smaller ε , the slower the convergence), or suffer from some numerical difficulties for small $\varepsilon > 0$.

Example 1. Let $d = 1$ and $\mu = \delta_0$, $\nu = \delta_\varepsilon$, $\varepsilon \neq 0$. For this toy example we know that $\|\mu - \nu\|_{TV} = 2$ and $(\phi_+^*, \phi_-^*) = (\mu, \nu)$. The semidefinite relaxation (2.9) with $n = 1$ reads:

$$\min_{\phi, \psi} \left\{ \phi_0 + \psi_0 : \phi_0 = \psi_0; \phi_1 - \psi_1 = -\varepsilon; \phi_2 - \psi_2 = -\varepsilon^2 \right. \\ \left. 0 \preceq \begin{bmatrix} \phi_0 & \phi_1 \\ \phi_1 & \phi_2 \end{bmatrix} \preceq \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}; 0 \preceq \begin{bmatrix} \psi_0 & \psi_1 \\ \psi_1 & -\psi_2 \end{bmatrix} \preceq \begin{bmatrix} 1 & \varepsilon \\ \varepsilon & \varepsilon^2 \end{bmatrix} \right\}.$$

The constraint $0 \preceq \mathbf{M}_1(\phi) \preceq \mathbf{M}_n(\delta_0)$ yields $1 \geq \phi_0$, and $\phi_2 = 0$, which in turn yields $\phi_1 = 0$. Hence $\psi_1 = a$ and $\psi_2 = a^2$. That is, the first semidefinite relaxation (2.9) provides the optimal solution (ϕ_+^*, ϕ_-^*) , no matter how close is ε to 0. We can see that crucial for the relaxations (2.9) are the domination constraints $\mathbf{M}_n(\phi) \preceq \mathbf{M}_n(\mu)$ and $\mathbf{M}_n(\psi) \preceq \mathbf{M}_n(\nu)$ (whereas they are not needed in the LP (2.5)).

3. CONCLUSION

We have provided a numerical scheme to approximate as closely as desired the total variation distance between two measures μ and ν on \mathbb{R}^d . To the best of the author's knowledge, this is the first systematic algorithmic procedure to address this problem under fairly general assumptions on μ and ν , as both to satisfy Carleman's condition or the easier to check sufficient condition (2.3). On the other hand, this procedure is still "ideal" as for convergence it requires to have access to all moments of μ and ν exactly, or at least a sufficiently large finite number of them to obtain a good approximation, which can be questionable in real applications .

REFERENCES

- [1] *Handbook on Semidefinite, Conic and Polynomial Optimization*, Internat. Ser. Oper. Res. Management Sci., vol 166, M. Anjos and J. B. Lasserre (eds.) Springer, New York, 2012.
- [2] A. Chambolle. An algorithm for total variation minimization and applications, *J. Math. Imag. Vision* 20, pp. 89–97, 2004.
- [3] T. F. Chan, G. H. Glob, P. Mulet. A nonlinear primal-dual method for total variation-based image restoration, *SIAM J. Sci. Comput.* 20 (6), pp. 1964–1977, 1999.
- [4] D. Henrion, M. Korda, J.B. Lasserre. *The Moment-SOS Hierarchy: Lectures in Probability, Statistics, Computational Geometry, Control and Nonlinear PDEs*, World Scientific, Singapore, 2020.
- [5] D. Henrion, J. B. Lasserre. Graph recovery from incomplete moment information, *Constr. Approx.* 56, pp. 165–187, 2022.
- [6] D. Henrion, J. B. Lasserre, J. Lofberg. Gloptipoly 3: moments, optimization and semidefinite programming, *Optim. Methods and Softwares* 24, pp. 761–779, 2009.
- [7] J. B. Lasserre. The Moment-SOS Hierarchy, in *Proceedings of the International Congress of Mathematicians (ICM 2018)*, vol 4, B. Sirakov, P. Ney de Souza and M. Viana (eds.), World Scientific, 2019, pp. 3773–3794.
- [8] J. B. Lasserre. *Moments, Positive Polynomials and Their Applications*, Imperial College Press, London, UK, 2009.
- [9] M. Markatou, Yang Chen. Non-quadratic distances in model assessment, *Entropy* 20 (6), 464, 2018.
- [10] G. Peyré, M. Cuturi. Computational Optimal Transport: With applications to Data Science, Found. Trends in Machine Learning 11 (5-6), pp. 355–607, 2019.
- [11] L. I. Rudin, S. Osher, E. Fatemi. Nonlinear total variation based noise removal algorithms, *Physica D* 60, pp. 259–268, 1992.
- [12] T. Weisser, B. Legat, C. Coey, L. Kapelevich, J.P. Vielma. Polynomial and moment optimization in Julia and Jump, Juliacon, 2019.

LAAS-CNRS AND TOULOUSE SCHOOL OF ECONOMICS (TSE), UNIVERSITY OF TOULOUSE,
 LAAS, 7 AVENUE DU COLONEL ROCHE, BP 54200, 31031 TOULOUSE CÉDEX 4, FRANCE
Email address: `lasserre@laas.fr`