



**HAL**  
open science

# Optimal quadrature-sparsification for integral operator approximation

Bertrand Gauthier, Johan a K Suykens

► **To cite this version:**

Bertrand Gauthier, Johan a K Suykens. Optimal quadrature-sparsification for integral operator approximation. 2016. hal-01416786v1

**HAL Id: hal-01416786**

**<https://hal.science/hal-01416786v1>**

Preprint submitted on 14 Dec 2016 (v1), last revised 20 Dec 2019 (v5)

**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.

# OPTIMAL QUADRATURE-SPARSIFICATION FOR INTEGRAL OPERATOR APPROXIMATION

BERTRAND GAUTHIER<sup>\*†</sup> AND JOHAN A.K. SUYKENS<sup>‡†</sup>

**Abstract.** We address the problem of designing sparse quadratures for the approximation of integral operators related to symmetric positive-semidefinite kernels. For a given kernel, we introduce the notion of squared-kernel discrepancy between two measures, consisting in the squared Hilbert-Schmidt norm of the difference between the integral operators defined from the underlying kernel and the two measures, the integral operators being viewed as operators on the underlying reproducing kernel Hilbert space. In the framework of integral operators defined from measures supported by a fixed set of points, sparsity of the approximate quadrature can be enforced through an  $\ell^1$ -type penalisation, and the computation of a penalised squared-kernel-discrepancy-optimal Nyström approximation is thus turned into a convex quadratic minimisation problem. The penalisation can be introduced under the form of a regularisation term or of a constraint, both formulations being equivalent; by analogy with spectral truncation, we in particular propose to penalise the trace of the approximate operator. The quadratic programs related to the regularised and constrained squared-kernel-discrepancy minimisation problems can be interpreted as the Lagrange dual formulations of distorted one-class support-vector machines related to the squared-kernel and the initial discrete measure. Numerical strategies for solving large-scale squared-kernel discrepancy minimisation problems are investigated and the efficiency of the approach is illustrated on a series of examples. We in particular demonstrate the ability of the proposed methodology to lead to accurate approximations of the main eigenpairs of kernel-matrices related to large-scale datasets.

**Key words.** sparse Nyström approximation, integral operator, RKHS, optimal quadrature-sparsification, squared-kernel discrepancy, L1-type penalisation, convex quadratic programming, one-class SVM.

**AMS subject classifications.** 47G10, 41A55, 46E22

**1. Introduction.** The eigendecomposition of an integral operator defined from a symmetric and positive-semidefinite kernel and a discrete measure supported by  $N$  points is numerically equivalent to the diagonalisation of a  $N \times N$  symmetric and positive-semidefinite matrix (notice that such a framework is encountered each time a pointwise quadrature is used to approximate an integral operator, or more simply when one aims at performing the eigendecomposition of a kernel-matrix). In the non-sparse case and for the direct “naive” approach, the amount of computations required to perform the diagonalisation of such a matrix scales as  $\mathcal{O}(N^3)$ , and therefore becomes quickly intractable when  $N$  increases; in addition, the memory required for the storage of the underlying matrix scales as  $\mathcal{O}(N^2)$ . By analogy with the classical framework of Nyström approximation (see, e.g., [8]), a common alternative consists in carrying out the eigendecomposition of a (weighted) principal submatrix of the initial matrix, with size  $n \ll N$ , and then expanding the result back up to dimension  $N$ ; in the integral-operator framework, this operation is related to the definition of a sparse approximation of the initial discrete measure (the support of the approximate measure being included in the support of the initial measure).

Naturally, the selection of the considered sparse measure has a strong impact on the quality of the induced approximation, raising questions relative to the selection of an “appropriate” measure. In our framework, since the initial discrete measure may already be interpreted as a pointwise quadrature, we refer to this problem as *optimal quadrature-sparsification*; we in particular consider the issue of computing a sparse measure leading to an accurate approximation of the main eigendirections of the initial integral operator. As a main feature, our approach is deterministic, relies on convex quadratic optimisation, and leads to approximations enjoying an optimality property. By dealing with integral operators and kernel-matrices, we can potentially handle large-scale problems (i.e., large values of  $N$ ) without suffering storage issues: we only require to have a quick access to any entry of the involved matrices, the required entries being computed on the fly.

To be more precise, we consider a general space  $\mathcal{X}$  and a symmetric and positive-semidefinite kernel  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ ; the underlying reproducing kernel Hilbert space (RKHS, see for instance [1]) of real-valued functions on  $\mathcal{X}$  is denoted by  $\mathcal{H}$ . For a set  $S = \{x_k\}_{k=1}^N$  of  $N$  points in  $\mathcal{X}$ , we denote by  $\mathbf{K}$  the underlying  $N \times N$  kernel-matrix, with  $i, j$ -th entry  $\mathbf{K}_{i,j} = K(x_i, x_j)$ . Considering a

<sup>\*</sup>bertrand.gauthier@esat.kuleuven.be (corresponding author)

<sup>‡</sup>johan.suykens@esat.kuleuven.be

<sup>†</sup>KU Leuven, ESAT-STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics  
Kasteelpark Arenberg 10, B-3001 Leuven, Belgium

measure  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  supported by  $\mathcal{S}$  (with  $\omega_k > 0$  and where  $\delta_{x_k}$  stands for Dirac measure centered at  $x_k$ ), we can naturally define the integral operator (see Section 2 for a detailed discussion)

$$T_\mu[f](x) = \int_{\mathcal{X}} K(x, t) f(t) d\mu(t).$$

The operator  $T_\mu$  is canonically related to the matrix  $\mathbf{KW}$  where  $\mathbf{W} = \text{diag}(\boldsymbol{\omega})$  is the  $N \times N$  diagonal matrix with diagonal entries  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_N)^T$ ; the eigendecomposition of  $T_\mu$  can be easily obtained from the diagonalisation of the symmetric positive-semidefinite matrix  $\mathbf{W}^{1/2} \mathbf{K} \mathbf{W}^{1/2}$  (notice that the diagonalisation of  $\mathbf{K}$  corresponds to the case  $\boldsymbol{\omega} = \mathbb{1}$ , with  $\mathbb{1} = (1, \dots, 1)^T \in \mathbb{R}^N$ ). Our aim is to construct a sparse measure  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$ , with  $v_k \geq 0$  ( $\nu$  is therefore also supported by  $\mathcal{S}$ ) in order to use the eigendecomposition of  $T_\nu$  to build an accurate approximation of the main eigenpairs of  $T_\mu$ . By “sparse measure  $\nu$ ”, we mean that we expect most of the components of the vector  $\mathbf{v} = (v_1, \dots, v_N)^T$  to be null, making the eigendecomposition of  $T_\nu$  easily computable; the measure  $\nu$  may therefore be interpreted as a sparse approximation of the initial measure  $\mu$ .

In order to construct the measure  $\nu$ , we propose to consider a constrained or regularised minimisation of the *squared-kernel discrepancy between  $\mu$  and  $\nu$* , denoted by  $D_{K^2}(\mu, \nu)$  and given by

$$D_{K^2}(\mu, \nu) = \|T_\mu - T_\nu\|_{\text{HS}, \mathcal{H}}^2,$$

i.e.,  $D_{K^2}(\mu, \nu)$  is the squared Hilbert-Schmidt norm of the difference between the operators  $T_\mu$  and  $T_\nu$  interpreted as integral operators on the RKHS  $\mathcal{H}$ , see Section 2.3 (an similar framework is for instance considered in [18]). The connection between the squared-kernel discrepancy and the Nyström approximation problem is further highlighted by Theorem 2.1.

When  $\mu$  and  $\nu$  are discrete measures supported by  $\mathcal{S}$  (related to the vectors  $\boldsymbol{\omega}$  and  $\mathbf{v} \in \mathbb{R}^N$ , respectively), we have  $D_{K^2}(\mu, \nu) = (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v})$ , where  $\mathbf{S} = \mathbf{K} * \mathbf{K}$  (Hadamard product) is the kernel-matrix defined by the (symmetric and positive-semidefinite) squared-kernel  $K^2(\cdot, \cdot) = (K(\cdot, \cdot))^2$  and  $\mathcal{S}$ . Remarkably, the squared-kernel discrepancy appears as a natural extension of the classical Frobenius-norm-based criterion for matrix low-rank approximation (see. e.g., [4, 11, 20, 2, 6]), as discussed in Section 3.1.

For a fixed measure  $\mu$ , the minimum of the application  $\nu \mapsto D_{K^2}(\mu, \nu)$  is trivially reached for  $\nu = \mu$ ; an  $\ell_1$ -type penalisation is therefore introduced in order to discard this trivial solution while enforcing the sparsity of the measure  $\nu$ , see Section 3. The penalisation takes the form of a regularisation term or of an additional constraint, both formulations being equivalent. In this way, the sampling problem for Nyström approximation is turned into a quadratic program (QP), and the sparsity of the resulting optimal measure  $\nu^*$  is controlled by the penalisation parameter (namely,  $\alpha$  for the regularised formulation, or  $\varkappa$  for the constrained formulation). By analogy with the framework of spectral truncation, we more particularly focus on the penalisation of the trace of the operator  $T_\nu$ , while for instance enforcing  $\text{trace}(T_\nu)$  to be a fraction of  $\text{trace}(T_\mu)$ .

Since the matrix  $\mathbf{S}$  is a kernel-matrix, the QPs related to the regularised or constrained discrepancy minimisation can naturally be interpreted as the Lagrange duals of “distorted” one-class support vector machines (SVMs) related to the squared-kernel  $K^2(\cdot, \cdot)$  and the initial discrete measured  $\mu$ , as discussed in Section 5. This analogy with one-class SVMs also suggests the introduction of soft-margin-type extensions of the squared-kernel-discrepancy minimisations problems.

We present a careful analysis of the approach and propose two different numerical strategies for the computations of optimal (or nearly optimal) sparse measures  $\nu^*$ . The first strategy, described in Section 6, is “direct” and consists in computing the regularisation path (see for instance [14]) related to the regularised or constrained formulations. The second strategy, see Section 7, applies more specifically to the constrained formulation and consists in using a vertex-exchange algorithm (see, e.g., [15, Chap. 9]) in order to solve the underlying QP. Importantly, in order to be able to tackle large-scale problems, our numerical strategies are “kernelised”, in the sense that they do not require the storage of any large objects (like in particular the matrix  $\mathbf{K}$  or  $\mathbf{S}$ ). In order to enhance the sparsity of a given discrete measure  $\nu$  while trying to keep the discrepancy as low as possible, in Section 7.2, we propose two greedy heuristics based on iterative pairwise-component merging.

Some experiments are carried out in Sections 8 and 9, illustrating the ability of the approach to lead to accurate sparse approximations of the main eigenpairs of any integral operator defined from a discrete measure with finite support.

The paper is organised as follows. Section 2 is devoted to some theoretical recalls about integral operators and to the introduction of the notion of squared-kernel discrepancy. In Section 3, we only consider discrete measures with support included in a fixed finite set  $S$ ; in this framework, we introduce the QPs related to the regularised and constrained squared-kernel-discrepancy minimisation problems and discuss their properties. The approximation of the main eigenpairs of the initial operator  $T_\mu$  from a low-discrepancy measure  $\nu$  is investigated in Section 4. Section 5 is devoted to the study of the one-class SVMs related to the regularised and constrained problems. Some numerical strategies for the resolution of the regularised and constrained problems are discussed in Sections 6 and 7. Sections 8 (two-dimensional toy example) and 9 (large-scale problems) are devoted to numerical experiments, and Section 10 concludes.

We have tried to make the paper as self-contained as possible; for the sake of readability, all the proofs are placed in Appendix A.

**2. Notations, recalls and theoretical motivations.** This section is devoted to recalls related to integral operators defined from a symmetric positive-semidefinite kernel; we then introduce the notion of squared-kernel discrepancy and highlight its relevance in the framework of Nyström approximation.

**2.1. Integral operators.** We assume  $\mathcal{X}$  is a measurable space and we denote by  $\mathcal{A}$  the underlying  $\sigma$ -algebra. We suppose that the (real-valued) kernel  $K(\cdot, \cdot)$  is measurable on  $\mathcal{X} \times \mathcal{X}$  for the product  $\sigma$ -algebra  $\mathcal{A} \otimes \mathcal{A}$  (see for instance [19, Chap. 4]), so that, in particular, the RKHS  $\mathcal{H}$  consists of measurable functions on  $\mathcal{X}$ . We also assume that the diagonal of  $K(\cdot, \cdot)$ , i.e., the function  $x \mapsto K(x, x)$  is measurable on  $(\mathcal{X}, \mathcal{A})$ . We denote by  $\mathcal{M}$  the set of all measures on  $(\mathcal{X}, \mathcal{A})$  and we introduce

$$\mathcal{T}(K) = \left\{ \mu \in \mathcal{M} \mid \int_{\mathcal{X}} K(x, x) d\mu(x) = \tau_\mu < +\infty \right\}.$$

For  $\mu \in \mathcal{T}(K)$ , we have  $K(\cdot, \cdot) \in L^2(\mu \otimes \mu)$  since in particular

$$\|K\|_{L^2(\mu \otimes \mu)}^2 = \int_{\mathcal{X} \times \mathcal{X}} (K(x, t))^2 d\mu(x) d\mu(t) \leq \tau_\mu^2.$$

In addition, for all  $h \in \mathcal{H}$ , we have  $h \in L^2(\mu)$  and  $\|h\|_{L^2(\mu)}^2 \leq \tau_\mu \|h\|_{\mathcal{H}}^2$ , i.e.,  $\mathcal{H}$  is continuously included in  $L^2(\mu)$ . We then define the symmetric and positive-semidefinite integral operator  $T_\mu$  on  $L^2(\mu)$ , given by, for  $f \in L^2(\mu)$  and  $x \in \mathcal{X}$ ,

$$T_\mu[f](x) = \int_{\mathcal{X}} K(x, t) f(t) d\mu(t),$$

so that  $T_\mu[f] \in \mathcal{H}$  and, for  $h \in \mathcal{H}$ ,  $(h|T_\mu[f])_{\mathcal{H}} = (h|f)_{L^2(\mu)}$ , see for instance [5] for more details.

We introduce the closed linear subspaces  $\mathcal{H}_{0\mu} = \{h \in \mathcal{H} \mid \|h\|_{L^2(\mu)} = 0\}$  and  $\mathcal{H}_\mu = \mathcal{H}_{0\mu}^{\perp}$  (i.e.,  $\mathcal{H}_\mu$  is the orthogonal of  $\mathcal{H}_{0\mu}$  in  $\mathcal{H}$ ), leading to the orthogonal decomposition  $\mathcal{H} = \mathcal{H}_\mu \oplus \mathcal{H}_{0\mu}$ .

We denote by  $\{\lambda_k \mid k \in \mathbb{N}^+\}$  the at most countable set of all strictly positive eigenvalues of  $T_\mu$  (repeated according to their algebraic multiplicity), and let  $\tilde{\varphi}_k \in L^2(\mu)$  be the associated eigenfunctions, orthonormalised for  $L^2(\mu)$ . For  $k \in \mathbb{N}^+$ , let  $\varphi_k = \frac{1}{\lambda_k} T_\mu[\tilde{\varphi}_k] \in \mathcal{H}$  be the canonical extension of  $\tilde{\varphi}_k$ , so that  $\{\sqrt{\lambda_k} \varphi_k \mid k \in \mathbb{N}^+\}$  is an orthonormal basis (o.n.b.) of the subspace  $\mathcal{H}_\mu$  of  $\mathcal{H}$ . The reproducing kernel of  $\mathcal{H}_\mu$  (for the Hilbert structure of  $\mathcal{H}$ ) is  $K_\mu(x, t) = \sum_{k \in \mathbb{N}^+} \lambda_k \varphi_k(x) \varphi_k(t)$  and  $\tau_\mu = \sum_{k \in \mathbb{N}^+} \lambda_k \geq 0$  (in fact,  $\tau_\mu > 0$  as far as  $\mathbb{N}^+$  is not empty).

**2.2. Discrete measures and kernel-matrices .** Let  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  be a discrete measure supported by  $S$ , with  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_N)^T \in \mathbb{R}^N$ ,  $\omega_k > 0$  (in what follows, we shall use the notation  $\boldsymbol{\omega} > 0$ ), where  $\delta_{x_k}$  is the Dirac measure (evaluation functional) at  $x_k$ . We have  $\mu \in \mathcal{T}(K)$ , and for  $f \in L^2(\mu)$  and  $x \in \mathcal{X}$ , using matrix notation,

$$T_\mu[f](x) = \sum_{k=1}^N \omega_k K(x, x_k) f(x_k) = \mathbf{k}^T(x) \mathbf{W} \mathbf{f},$$

with  $\mathbf{W} = \text{diag}(\boldsymbol{\omega})$ ,  $\mathbf{k}(x) = (K(x_1, x), \dots, K(x_N, x))^T$  and  $\mathbf{f} = (f(x_1), \dots, f(x_N))^T \in \mathbb{R}^N$ . We can identify the Hilbert space  $L^2(\mu)$  with the space  $\mathbb{R}^N$  endowed with the inner product  $(\cdot|\cdot)_{\mathbf{W}}$ , where for  $\mathbf{x}$  and  $\mathbf{y} \in \mathbb{R}^N$ ,  $(\mathbf{x}|\mathbf{y})_{\mathbf{W}} = \mathbf{x}^T \mathbf{W} \mathbf{y}$ ; in this way,  $f \in L^2(\mu)$  is assimilated to  $\mathbf{f} \in \mathbb{R}^N$ , and the operator  $T_\mu$  then corresponds to the matrix  $\mathbf{K} \mathbf{W}$ . We denote by  $\lambda_1 \geq \dots \geq \lambda_N \geq 0$  the eigenvalues of  $\mathbf{K} \mathbf{W}$  and by  $\mathbf{v}_1, \dots, \mathbf{v}_N$  the associated orthonormalised eigenvectors, i.e.,  $\mathbf{K} \mathbf{W} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^{-1}$  with  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$  and  $\mathbf{P} = (\mathbf{v}_1 | \dots | \mathbf{v}_N)$ . The vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_N\}$  form an orthonormal basis of the Hilbert space  $\{\mathbb{R}^N, (\cdot|\cdot)_{\mathbf{W}}\}$ , i.e.,  $\mathbf{P}^T \mathbf{W} \mathbf{P} = \text{Id}_N$ , the  $N \times N$  identity matrix. In particular, we have  $\mathbf{K} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$ , and  $\mathbf{P}^T \mathbf{P} = \mathbf{W}^{-1}$ . For  $\lambda_k > 0$ , the canonical eigenfunctions of  $T_\mu$  are given by  $\varphi_k = \frac{1}{\lambda_k} \mathbf{k}^T \mathbf{W} \mathbf{v}_k$ .

For a general  $\boldsymbol{\omega} > 0$ , the matrix  $\mathbf{K} \mathbf{W}$  is non-symmetric and computing its diagonalisation may be numerically difficult. However, classically, since  $\mathbf{K} \mathbf{W} \mathbf{v}_k = \lambda_k \mathbf{v}_k$ , we have

$$\mathbf{W}^{1/2} \mathbf{K} \mathbf{W}^{1/2} \mathbf{W}^{1/2} \mathbf{v}_k = \lambda_k \mathbf{W}^{1/2} \mathbf{v}_k.$$

The symmetric matrix  $\mathbf{W}^{1/2} \mathbf{K} \mathbf{W}^{1/2}$  thus defines a symmetric and positive-semidefinite operator on  $\{\mathbb{R}^N, (\cdot|\cdot)_{\text{Id}_N}\}$  with eigenvalues  $\lambda_k$  and orthonormalised eigenvectors  $\mathbf{W}^{1/2} \mathbf{v}_k$ . One can therefore easily deduce the eigendecomposition of the matrix  $\mathbf{K} \mathbf{W}$  viewed as an operator on  $\{\mathbb{R}^N, (\cdot|\cdot)_{\mathbf{W}}\}$  from the eigendecomposition the symmetric matrix  $\mathbf{W}^{1/2} \mathbf{K} \mathbf{W}^{1/2}$ .

**2.3. Hilbert-Schmidt norm and squared-kernel discrepancy.** In view of Section 2.1, for  $\mu \in \mathcal{T}(K)$ , the operator  $T_\mu$  can also be interpreted as an operator on  $\mathcal{H}$ ; with a slight abuse of notation, we keep the same notation for “ $T_\mu$  viewed as an operator on  $L^2(\mu)$ ”, and “ $T_\mu$  viewed as an operator on  $\mathcal{H}$ ”. In both cases,  $T_\mu$  is an Hilbert-Schmidt operator.

Let  $\mu$  and  $\nu \in \mathcal{T}(K)$ ; for an o.n.b.  $\{h_k | k \in \mathbb{I}\}$  of  $\mathcal{H}$  (with  $\mathbb{I}$  a general, at most countable, index set), the Hilbert-Schmidt inner product between the operators  $T_\mu$  and  $T_\nu$  on  $\mathcal{H}$  is given by

$$(T_\mu | T_\nu)_{\text{HS}, \mathcal{H}} = \sum_{k \in \mathbb{I}} (T_\mu[h_k] | T_\nu[h_k])_{\mathcal{H}},$$

and we recall that the value of  $(T_\mu | T_\nu)_{\text{HS}, \mathcal{H}}$  does not depend on the choice of the o.n.b. of  $\mathcal{H}$ , see, e.g., [17]. The underlying Hilbert-Schmidt norm (for operators on  $\mathcal{H}$ ) is given by

$$\|T_\mu\|_{\text{HS}, \mathcal{H}}^2 = \sum_{k \in \mathbb{I}} \|T_\mu[h_k]\|_{\mathcal{H}}^2.$$

**Lemma 2.1.** *For  $\mu$  and  $\nu \in \mathcal{T}(K)$ , we have  $(T_\mu | T_\nu)_{\text{HS}, \mathcal{H}} = \|K\|_{L^2(\mu \otimes \nu)}^2$ . In addition, we have  $\|T_\mu\|_{\text{HS}, \mathcal{H}}^2 = \sum_{k \in \mathbb{I}_\mu^+} \lambda_k^2$ , where  $\{\lambda_k | k \in \mathbb{I}_\mu^+\}$  is the set of all strictly positive eigenvalues of  $T_\mu$ .*

By definition, we also remark that if  $\mu$  and  $\nu \in \mathcal{T}(K)$  are such that  $\mathcal{H}_\mu$  and  $\mathcal{H}_\nu$  are orthogonal subspaces of  $\mathcal{H}$ , then  $\|K\|_{L^2(\mu \otimes \nu)}^2 = 0$ .

**Definition 2.1.** *For  $\mu$  and  $\nu \in \mathcal{T}(K)$ , the squared-kernel discrepancy between  $\mu$  and  $\nu$  is*

$$D_{K^2}(\mu, \nu) = \|T_\mu - T_\nu\|_{\text{HS}, \mathcal{H}}^2 = \|K\|_{L^2(\mu \otimes \mu)}^2 + \|K\|_{L^2(\nu \otimes \nu)}^2 - 2\|K\|_{L^2(\mu \otimes \nu)}^2.$$

The terminology “discrepancy” is motivated by the analogy with the notion of “kernel discrepancy” used for instance in [3]. By definition,  $D_{K^2}(\mu, \nu)$  appears as a “measure of the similarity” between the operators  $T_\mu$  and  $T_\nu$  on  $\mathcal{H}$ ; in particular,  $D_{K^2}(\mu, \mu) = 0$ . Notice that  $\|K\|_{L^2(\mu \otimes \nu)}^2 \leq \tau_\mu \tau_\nu$ .

**2.4. Squared-kernel discrepancy and Nyström approximation error.** The following Theorem 2.1 further highlights the connection between the squared-kernel discrepancy and the error induced by the approximation of  $T_\mu$  (the “true” integral operator) by  $T_\nu$  (the approximate operator).

**Theorem 2.1.** *Let  $\mu$  and  $\nu \in \mathcal{T}(K)$  be such that  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$  (i.e., for  $h \in \mathcal{H}$ , if  $\|h\|_{L^2(\mu)} = 0$  then  $\|h\|_{L^2(\nu)} = 0$ ). Let  $\{\sqrt{\lambda_k} \varphi_k | k \in \mathbb{I}_\mu^+\}$  be the o.n.b. basis of  $\mathcal{H}_\mu$  defined by  $T_\mu$ , and  $\{\sqrt{\vartheta_l} \psi_l | l \in \mathbb{I}_\nu^+\}$  be the o.n.b. of  $\mathcal{H}_\nu$  defined by  $T_\nu$ . We have*

$$D_{K^2}(\mu, \nu) = \sum_{k \in \mathbb{I}_\mu^+} \lambda_k \|T_\mu[\varphi_k] - T_\nu[\varphi_k]\|_{\mathcal{H}}^2 \quad (2.1)$$

$$= \sum_{k \in \mathbb{I}_\mu^+} \lambda_k^2 + \sum_{k \in \mathbb{I}_\mu^+} \sum_{l \in \mathbb{I}_\nu^+} [\vartheta_l^2 - 2\lambda_k \vartheta_l] (\sqrt{\lambda_k} \varphi_k | \sqrt{\vartheta_l} \psi_l)_{\mathcal{H}}^2, \quad (2.2)$$

and, in addition,  $\sum_{k \in \mathbb{I}_\mu^+} \lambda_k \|T_\mu[\varphi_k] - T_\nu[\varphi_k]\|_{L^2(\mu)}^2 \leq \tau_\mu D_{K^2}(\mu, \nu)$ .

For a fixed measure  $\mu$ , by minimising the function  $\nu \mapsto D_{K^2}(\mu, \nu)$  under the constraint  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$ , we aim at minimising, for the RKHS norm,  $\lambda_k \varphi_k - T_\nu[\varphi_k]$  for all  $k \in \mathbb{I}_\mu^+$ , while putting more emphasis on the approximation of the eigenpairs with largest eigenvalues (the eigenvalues  $\lambda_k$  playing the role of penalisation weights). We can therefore in particular expect the main eigenfunctions of  $T_\nu$  to be accurate approximations of the main eigendirections of  $T_\mu$  (i.e., the one related to the largest eigenvalues). However, under the only condition  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$ ,  $D_{K^2}(\mu, \nu)$  is minimal for  $\nu = \mu$  (meaning that the best approximation of  $T_\mu$  is  $T_\mu$  itself); in order to avoid this trivial solution, one may add additional feasibility constraints for  $\nu$ , or add a regularisation term to the cost  $D_{K^2}(\mu, \nu)$ . By analogy with spectral truncation, one possibility is to penalise  $\tau_\nu = \text{trace}(T_\nu) = \int_{\mathcal{X}} K(x, x) d\nu(x)$ , see Section 3; more generally, in this work, we consider penalisations of the form  $\int_{\mathcal{X}} f(x) d\nu(x)$ , where  $f$  is a real-valued, positive, measurable function on  $\mathcal{X}$ .

**3. Optimal quadrature-sparsification as quadratic programming.** We assume that the measure  $\mu$  is discrete and has a finite number of support points (framework of Section 2.2), and we enforce the support of the searched after approximation  $\nu$  to be included in the support of  $\mu$ . As already mentioned, since  $D_{K^2}(\mu, \mu) = 0$ , the minimisation of the squared-kernel discrepancy alone is not interesting (i.e., minimisation of  $\nu \mapsto D_{K^2}(\mu, \nu)$  for  $\mu$  fixed); in addition, we aim at obtaining as sparse as possible discrete measure  $\nu$ . In order to discard the trivial solution  $\nu = \mu$ , we introduce a linear penalisation.

More precisely, let  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  and  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$  be discrete measures supported by  $\mathcal{S}$ , with  $\boldsymbol{\omega} > 0$  and  $\mathbf{v} \geq 0$  (the notation  $\mathbf{v} \geq 0$  meaning that all the components of  $\mathbf{v}$  are positive); we thus have  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$  for any such  $\mathbf{v}$  and

$$D_{K^2}(\mu, \nu) = (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v}),$$

where  $\mathbf{S}$  is the kernel-matrix defined by the squared-kernel  $K^2(\cdot, \cdot)$  and the set of points  $\mathcal{S}$ , i.e., with  $i, j$  entry  $K^2(x_i, x_j) \geq 0$  ( $\mathbf{S}$  is therefore a non-negative, positive-semidefinite, symmetric matrix). For a fixed measure  $\mu$  (i.e.,  $\boldsymbol{\omega}$  is fixed), we define

$$D(\mathbf{v}) = \frac{1}{2} (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v}),$$

the factor  $1/2$  being added for simplification purpose.

**3.1. Relation with the classical Frobenius-norm-based criterion.** In the framework of Section 3, the squared-kernel discrepancy appears as a natural extension of the classical Frobenius-norm-based criterion for matrix low-rank approximation (see for instance [4]). Indeed, assume that  $\boldsymbol{\omega} = \mathbf{1}$ , with  $\mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^N$ , so that  $\mathbf{W} = \text{Id}_N$  is the  $N \times N$  identity matrix, and thus  $\mathbf{KW} = \mathbf{W}^{1/2} \mathbf{KW}^{1/2} = \mathbf{K}$  (notice the analogy with the notations used in Section 2.2).

Selecting columns of the matrix  $\mathbf{K}$  can be viewed as performing the product  $\mathbf{KV}$ , where  $\mathbf{V} = \text{diag}(\mathbf{v})$  is a  $N \times N$  diagonal matrix with diagonal entries 0 or 1 (we are therefore also assuming that all the components of  $\mathbf{v}$  are 0 or 1); in the same way, considering  $\mathbf{VK}$  amounts to selecting rows of  $\mathbf{K}$ . Since such a sampling matrix  $\mathbf{V}$  satisfies  $\mathbf{V}^2 = \mathbf{V}$ , selecting a principal submatrix of the matrix  $\mathbf{K}$  can be written  $\mathbf{V}^{1/2} \mathbf{KV}^{1/2} = \mathbf{VKV}$ .

**Lemma 3.1.** *If  $\boldsymbol{\omega} = \mathbf{1}$  and if the components of  $\mathbf{v}$  are all 0 or 1, then  $(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v}) = \|\mathbf{K} - \mathbf{VKV}\|_F^2$ , where  $\|\cdot\|_F$  stands for the Frobenius norm.*

**3.2. Regularised squared-kernel-discrepancy minimisation.** For a given penalisation vector  $\mathbf{d} = (d_1, \dots, d_N)^T \in \mathbb{R}^N$ ,  $\mathbf{d} > 0$  (see Section 3.4), and  $\alpha \geq 0$ , we introduce the regularised squared-kernel-discrepancy minimisation problem

$$\underset{\mathbf{v}}{\text{minimise}} \quad D_\alpha(\mathbf{v}) = \frac{1}{2} (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v}) + \alpha \mathbf{d}^T \mathbf{v} \quad \text{subject to } \mathbf{v} \geq 0. \quad (3.1)$$

Notice that  $D_\alpha(\mathbf{v}) = D(\mathbf{v}) + \alpha \mathbf{d}^T \mathbf{v}$ . In particular, when  $\mathbf{S}$  is invertible,  $D_\alpha(\cdot)$  is strongly convex (and, in this case, a solution to (3.1) is therefore necessarily unique). We also recall that, for a given  $\alpha$ , the set of solutions to (3.1) is convex. The gradient of  $D_\alpha$  at  $\mathbf{v} \in \mathbb{R}^N$  is given by  $\nabla D_\alpha(\mathbf{v}) = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega}) + \alpha \mathbf{d}$ . The following Lemma 3.2 recalls some simple properties verified by the solutions to problem (3.1).

**Lemma 3.2.** Denote by  $\mathbf{v}_\alpha^*$  a solution to (3.1) with  $\alpha \geq 0$ , we have:

- (i) for  $\alpha = 0$ ,  $\mathbf{v}_\alpha^* = \boldsymbol{\omega}$  is a solution to (3.1),
- (ii) if  $\alpha \geq \max_k \{[\mathbf{S}\boldsymbol{\omega}]_k/d_k\}$ , then  $\mathbf{v}_\alpha^* = 0$  (with  $[\mathbf{S}\boldsymbol{\omega}]_k$  the  $k$ -th component of  $\mathbf{S}\boldsymbol{\omega}$ ),
- (iii) for all  $\alpha \geq 0$ , we have  $0 \leq \alpha \mathbf{d}^T \mathbf{v}_\alpha^* \leq \alpha \mathbf{d}^T \boldsymbol{\omega} - (\boldsymbol{\omega} - \mathbf{v}_\alpha^*)^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}_\alpha^*)$ ,
- (iv) if  $\tilde{\mathbf{v}}_\alpha^*$  is another solution to (3.1), then  $\mathbf{d}^T \tilde{\mathbf{v}}_\alpha^* = \mathbf{d}^T \mathbf{v}_\alpha^*$ ,
- (v) if  $[\alpha \mathbf{d} - \mathbf{S}\boldsymbol{\omega}]_k \geq 0$  and  $\mathbf{S}_{k,k} > 0$  (see Remark 3.1), then  $[\mathbf{v}_\alpha^*]_k = 0$ ,
- (vi) the map  $\alpha \mapsto D(\mathbf{v}_\alpha^*)$  is increasing, and  $\alpha \mapsto \mathbf{d}^T \mathbf{v}_\alpha^*$  is decreasing.

Since  $\mathbf{v} \geq 0$ , the term  $\mathbf{d}^T \mathbf{v}$  can be interpreted as a weighted  $\ell^1$ -regularisation (and  $\alpha$  is the regularisation parameter). For appropriate  $\mathbf{d}$  and  $\alpha$ , we can thus expect a solution  $\mathbf{v}_\alpha^*$  to (3.1) to be sparse (see, e.g., [10]). This intuition is confirmed by Lemma 3.2-(v); indeed, when all the diagonal entries of  $\mathbf{S}$  are strictly positive,  $\text{card}(\{k | [\alpha \mathbf{d} - \mathbf{S}\boldsymbol{\omega}]_k < 0\})$  gives an upper bound on the number of strictly positive components of  $\mathbf{v}_\alpha^*$  (notice that this bound is generally not tight). See Section 3.4 for a further discussion.

**Remark 3.1.** Assuming  $\mathbf{S}_{k,k} = K^2(x_k, x_k) > 0$  for all  $k \in \{1, \dots, N\}$  (which we denote by  $\text{diag}(\mathbf{S}) > 0$ ) is equivalent to assuming  $K(x_k, x_k) > 0$  for all  $k$ ; we recall that for all  $x \in \mathcal{X}$ , we have  $K(x, x) = \|K_x\|_{\mathcal{H}}^2 \geq 0$ . This assumption is in practice not restrictive at all: indeed, if  $K(x_k, x_k) = 0$ , then  $K_{x_k} = 0$  and thus  $h(x_k) = 0$  for all  $h \in \mathcal{H}$ . In the framework of Section 3, such a point  $x_k$  may thus be removed from the sample  $\mathcal{S}$  without inducing any modification of the operators  $T_\mu$  and  $T_\nu$ .  $\triangleleft$

**3.3. Constrained squared-kernel-discrepancy minimisation.** Instead of considering problem (3.1), we can equivalently introduce, for  $\varkappa \geq 0$  (and  $\varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$ , see Lemma 3.3)

$$\text{minimise } D(\mathbf{v}) = \frac{1}{2}(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) \text{ subject to } \mathbf{v} \geq 0 \text{ and } \mathbf{d}^T \mathbf{v} = \varkappa. \quad (3.2)$$

Notice that problem (3.2) consists in minimising a convex function on a convex compact domain, so that a solution  $\mathbf{v}_\varkappa^*$  to problem (3.2) always exists; in particular, if  $\mathbf{S}$  is non-singular, then  $\mathbf{v}_\varkappa^*$  is always unique (since  $D(\cdot)$  is in this case strongly convex).

**Lemma 3.3.** Let  $\mathbf{v}_\alpha^*$  be a solution to problem (3.1) with  $\alpha \geq 0$ ; then  $\mathbf{v}_\alpha^*$  is a solution to problem (3.2) with  $\varkappa = \mathbf{d}^T \mathbf{v}_\alpha^*$ . Reciprocally, assume that  $\mathbf{v}_\varkappa^*$  is a solution to problem (3.2) with  $0 < \varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$ ; then  $\mathbf{v}_\varkappa^*$  is a solution to problem (3.1) with  $\alpha = (\mathbf{v}_\varkappa^*)^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}_\varkappa^*)/\varkappa$ . For  $\varkappa = 0$ , we have  $\mathbf{v}_\varkappa^* = 0$ , which is solution to problem (3.1) with  $\alpha \geq \max_k \{[\mathbf{S}\boldsymbol{\omega}]_k/d_k\}$ . For  $0 \leq \varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$ , the map  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  is decreasing.

As an interesting feature, problem (3.2) can be efficiently solved thanks to a kernelised (i.e., without matrix storage) sparse descent direction QP solver, like for instance the vertex-exchange strategy, see [15, Chap. 9] and Section 7.1. A sequential strategy (based on the notion of regularisation path) for solving problems (3.1) and (3.2) is also discussed in Section 6

Notice that, in view of Lemma 3.2-(iii) and Lemma 3.3, considering  $\varkappa = \rho \mathbf{d}^T \boldsymbol{\omega}$  with  $\rho \in [0, 1]$  appears as a natural parameterisation for problem (3.2).

**Remark 3.2.** In case of non-uniqueness of the solution to problem (3.2) with  $\varkappa \geq 0$ , two solutions  $\mathbf{v}_\varkappa^*$  and  $\tilde{\mathbf{v}}_\varkappa^*$  necessarily satisfy  $\mathbf{S}(\mathbf{v}_\varkappa^* - \tilde{\mathbf{v}}_\varkappa^*) = 0$ . In exactly the same way, from Lemma 3.2-(iv) or Lemma 3.3, in case of non-uniqueness of the solution to problem (3.1) with  $\alpha \geq 0$ , two solutions  $\mathbf{v}_\alpha^*$  and  $\tilde{\mathbf{v}}_\alpha^*$  necessarily satisfy  $\mathbf{S}(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*) = 0$ .  $\triangleleft$

**3.4. Choice of the penalisation.** In order to discard the solution  $\mathbf{v} = \boldsymbol{\omega}$  (i.e.,  $\nu = \mu$ ), in Sections 3.2 and 3.3, we have introduced a penalisation based on the term  $\mathbf{d}^T \mathbf{v}$  (the penalisation appearing in the cost or in the constraint), for a given  $\mathbf{d} \in \mathbb{R}^N$  with  $\mathbf{d} > 0$ . Notice that the term  $\mathbf{d}^T \mathbf{v}$  can be interpreted as the integral with respect to  $\nu$  of a real valued function  $f_{\mathbf{d}}$  on  $\mathcal{X}$  satisfying  $f_{\mathbf{d}}(x_k) = d_k$ ; in this case, we indeed have  $\mathbf{d}^T \mathbf{v} = \int_{\mathcal{X}} f_{\mathbf{d}}(x) d\nu(x)$ . In practice, we ideally aim at obtaining a vector  $\mathbf{v}^*$  which is both as sparse as possible and such that  $D(\mathbf{v}^*)$  is as small as possible; this naturally raises questions relative to the choice of the penalisation vector  $\mathbf{d}$  (see also Section 10).

We can first remark that for  $\mathbf{d} = \theta \mathbf{S}\boldsymbol{\omega}$  with  $\theta > 0$  and for  $\alpha \leq 1/\theta$ ,  $\mathbf{v}_\alpha^* = (1 - \alpha\theta)\boldsymbol{\omega}$  is a solution to (3.1) (and  $\mathbf{v}_\alpha^* = 0$  for  $\alpha > 1/\theta$ ); the solution  $\mathbf{v}_\alpha^*$  is therefore non-sparse, and such a choice for the penalisation vector  $\mathbf{d}$  appears to be of little interest. More generally, for  $\mathbf{d} = \mathbf{S}\boldsymbol{\eta} > 0$  (with  $\boldsymbol{\eta} \in \mathbb{R}^N$ ), a similar remark holds for all  $\alpha$  such that  $\boldsymbol{\omega} - \alpha\boldsymbol{\eta} \geq 0$ ; indeed, we have in this case  $\mathbf{v}_\alpha^* = \boldsymbol{\omega} - \alpha\boldsymbol{\eta}$ .

As mentioned in Remark 3.1, we may reasonably assume that  $K(x_k, x_k) > 0$  for all  $k$ . In this case, by analogy with spectral truncation, choosing  $d_k = K(x_k, x_k)$  appears as relatively interesting since it leads to  $\mathbf{d}^T \mathbf{v} = \text{trace}(T_v)$  (and  $\text{trace}(T_\mu) = \mathbf{d}^T \boldsymbol{\omega}$ ); in this case, we shall use the notation  $\mathbf{d} = \text{diag}(\mathbf{K})$ . This choice therefore appears as relatively natural and interesting and is the one considered in the experiments of Sections 8 and 9.

**4. Approximate eigendecomposition.** We consider the framework of Section 3. Let  $\mathbf{v} \geq 0$  with  $\mathbf{v} \neq 0$  (we recall that  $\mathbf{v} \in \mathbb{R}^N$ ); in practice,  $\mathbf{v}$  will be a solution (or an approximate solution) to problem (3.1) or (3.2) (or of one of the extensions discussed in Sections 5.3), with  $\mathbf{d} > 0$  and  $\alpha \geq 0$  or  $0 \leq \varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$  (more particularly, since we are not interested in the case  $\mathbf{v} = 0$ , we shall consider  $\alpha < \max_k [\mathbf{S}\boldsymbol{\omega}]_k / d_k$  or  $\varkappa > 0$ ). We recall that we denote by  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$  the measure related to  $\mathbf{v}$ .

We introduce the index set  $I = \{i | v_i > 0\}$  and let  $n = \text{card}(I)$  be the number of strictly positive components of  $\mathbf{v}$ ; we have in particular  $\nu = \sum_{i \in I} v_i \delta_{x_i}$  (i.e., we have discarded all the points  $x_k$  such that  $v_k = 0$ ). Following Section 2.4, the strictly positive eigenvalues  $\{\vartheta_k | k \in \mathbb{I}_v^+\}$  of  $T_v$  and their associated canonically extended eigenfunction  $\psi_k \in \mathcal{H}$ , orthonormalised for  $L^2(\nu)$ , can be easily obtained from the eigendecomposition of the  $n \times n$  (symmetric and positive-semidefinite) principal submatrix  $[\mathbf{V}^{1/2} \mathbf{K} \mathbf{V}^{1/2}]_{I,I}$ . We thus obtain the o.n.b.  $\{\sqrt{\vartheta_l} \psi_l | l \in \mathbb{I}_v^+\}$  of the subspace  $\mathcal{H}_v$  of  $\mathcal{H}$  related to  $T_v$ .

Depending on the considered application, (i.e., depending on the reason why the eigendecomposition of  $T_\mu$  is needed), we may directly use the eigenpairs  $\{(\vartheta_l, \psi_l) | l \in \mathbb{I}_v^+\}$  related to  $T_v$  in order to approximate the eigenpairs  $\{(\lambda_k, \varphi_k) | k \in \mathbb{I}_\mu^+\}$  related to  $T_\mu$ . However, it is also possible to define more relevant approximations. In what follows, we discuss possible ways to approximate the main eigenpairs of the integral operator  $T_\mu$  from the eigendecomposition of the integral operator  $T_v$ .

*Normalised approximated eigenfunctions.* Since the eigenfunctions  $\varphi_k$  of  $T_\mu$ , with  $k \in \mathbb{I}_\mu^+$ , are such that  $\|\widehat{\varphi}_l\|_{L^2(\mu)} = 1$ , for all  $l \in \mathbb{I}_v^+$ , we can normalise the eigenfunctions  $\psi_k$  as well. We thus introduce the *normalised approximated eigenfunctions* of  $T_\mu$  induced by  $T_v$ , given by, for  $l \in \mathbb{I}_v^+$ ,

$$\widehat{\varphi}_l = \psi_l / \|\psi_l\|_{L^2(\mu)}, \quad (4.1)$$

so that  $\|\widehat{\varphi}_l\|_{L^2(\mu)} = 1$ . Notice however that the true eigenfunctions  $\{\varphi_k | k \in \mathbb{I}_\mu^+\}$  are orthogonal in  $L^2(\mu)$ , while the approximated eigenfunctions  $\{\widehat{\varphi}_l | l \in \mathbb{I}_v^+\}$  are orthogonal in  $L^2(\nu)$ .

**Remark 4.1.** Controlling the orthogonality, in  $L^2(\mu)$ , between the approximate eigenfunctions  $\widehat{\varphi}_l$  appears as a relatively inexpensive way to evaluate the accuracy of the approximate eigenfunctions. Indeed, accurate approximate eigenfunctions  $\widehat{\varphi}_l$  should be almost mutually orthogonal in  $L^2(\mu)$ ; this condition is however only a necessary condition.  $\triangleleft$

*Rescaled eigenvalues.* Let  $\rho \geq 0$  be such that  $\text{trace}(T_v) = \rho \text{trace}(T_\mu)$ ; if  $\rho > 0$  (i.e.,  $\mathbf{v} \neq 0$  and  $\text{diag}(\mathbf{K}) > 0$ ), for  $l \in \mathbb{I}_v^+$ , we may introduce the rescaled eigenvalues  $\vartheta_l / \rho$ , so that  $\sum_{l \in \mathbb{I}_v^+} \vartheta_l / \rho = \text{trace}(T_\mu)$ . As an interesting feature, such a rescaling involves any expensive computation, but the quality of the induced approximation is in general relatively poor.

One may also consider a corrected rescaling factor  $\rho_c$ , taking for instance into account that we approximate only the  $n_v$  largest eigenvalues of  $T_\mu$ , where  $n_v = \text{card}(\mathbb{I}_v^+)$  is the number of strictly positive eigenvalues of  $T_v$ . Depending on the application, another possibility may consist in rescaling the spectrum in such a way that the largest rescaled eigenvalue is equal to a given value.

*Upsilon-test and induced approximated eigenvalues.* By definition, for all  $k \in \mathbb{I}_\mu^+$ , we have

$$(\varphi_k | T_\mu[\varphi_k])_{L^2(\mu)} / \|T_\mu[\varphi_k]\|_{L^2(\mu)} = 1 \text{ and } \lambda_l = \|T_\mu[\varphi_k]\|_{L^2(\mu)}.$$

By analogy, for all  $l \in \mathbb{I}_v^+$ , considering the normalised approximate eigenfunctions  $\widehat{\varphi}_l$  given by (4.1), we define

$$\Upsilon_l = (\widehat{\varphi}_l | T_\mu[\widehat{\varphi}_l])_{L^2(\mu)} / \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}. \quad (4.2)$$

Since  $T_\mu$  is a positive-semidefinite operator on  $L^2(\mu)$ , we have  $\Upsilon_l \geq 0$ ; in addition, from the Cauchy-Schwarz inequality, we have  $\Upsilon_l \leq 1$ , so that  $\Upsilon_l \in [0, 1]$ . Notice that  $\Upsilon_l$  simply consists in the  $L^2(\mu)$



inner product between the two normalised vector  $\widehat{\varphi}_l$  and  $T_\mu[\widehat{\varphi}_l]/\|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}$ . We also introduce

$$\widehat{\lambda}_l = \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}. \quad (4.3)$$

So, the closer  $Y_l$  is to 1, the closer  $\widehat{\varphi}_l$  is to be an eigenfunction of  $T_\mu$ . In addition, for  $Y_l \approx 1$ ,  $\widehat{\lambda}_l$  appears as a relevant approximation of the underlying eigenvalue of  $T_\mu$  (i.e., the eigenvalue related to the eigenfunction approximated by  $\widehat{\varphi}_l$ ). This intuition is confirmed by the following Lemma 4.1.

**Lemma 4.1.** *For  $l \in \mathbb{I}_v^+$ , consider  $\widehat{\varphi}_l$ ,  $Y_l$  and  $\widehat{\lambda}_l$  given by (4.1), (4.2) and (4.3) respectively; we have*

$$\|T_\mu[\widehat{\varphi}_l] - \widehat{\lambda}_l \widehat{\varphi}_l\|_{L^2(\mu)}^2 = 2\widehat{\lambda}_l^2(1 - Y_l), \text{ and } \|T_\mu[\widehat{\varphi}_l] - \widehat{\lambda}_l \widehat{\varphi}_l\|_{\mathcal{H}}^2 = \widehat{\lambda}_l(Y_l + \|\widehat{\lambda}_l^{1/2} \widehat{\varphi}_l\|_{\mathcal{H}}^2 - 2).$$

From Lemma 4.1, since  $\|T_\mu[\widehat{\varphi}_l] - \widehat{\lambda}_l \widehat{\varphi}_l\|_{\mathcal{H}}^2 \geq 0$ , we can also deduce that for  $l \in \mathbb{I}_v^+$ , we have  $\|\widehat{\lambda}_l^{1/2} \widehat{\varphi}_l\|_{\mathcal{H}}^2 \geq 2 - Y_l$ ; we recall that by definition,  $\|\sqrt{\lambda_k} \varphi_k\|_{\mathcal{H}}^2 = 1$ , for all  $k \in \mathbb{I}_\mu^+$ .

The Upsilon-test appears has a very interesting and simple tool since it indicates the accuracy of the approximate eigenpairs of  $T_\mu$  induced by  $T_\nu$ . However, for  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$ , computing  $T_\mu[\widehat{\varphi}_l]$  consists in a kernelised matrix-vector product, with complexity scaling as  $\mathcal{O}(N^2)$ ; this operation is therefore relatively costly, but may however be easily parallelised.

*RKHS-norm-induced approximate eigenvalues.* As already mentioned, the numerical complexity of the Upsilon-test (and therefore, of the computation of the approximate eigenvalues  $\widehat{\lambda}_l$ ) may become prohibitive for very large  $N$  (depending on the computational power at disposal). One can however also define approximate eigenvalues involving computations with numerical complexity scaling as  $\mathcal{O}(n^2)$ , with  $n$  the number of support points of the considered approximate measure  $\nu$ .

Indeed, since  $\|\sqrt{\lambda_k} \varphi_k\|_{\mathcal{H}}^2 = 1$ , we have  $\lambda_k = 1/\|\varphi_k\|_{\mathcal{H}}^2$ . This therefore suggests to define the approximate eigenvalues, for  $l \in \mathbb{I}_v^+$ ,

$$\widetilde{\lambda}_l = 1/\|\widehat{\varphi}_l\|_{\mathcal{H}}^2. \quad (4.4)$$

From Section 2.2, the computation of the approximate eigenvalues  $\widetilde{\lambda}_l$  only involves the  $n \times n$  principal kernel-submatrix  $\mathbf{K}_{l,l}$ , and may in addition be easily parallelised. If  $\widehat{\varphi}_l$  is close to be a true normalised eigenfunction of  $T_\mu$ , then we can expect  $\widetilde{\lambda}_l$  to be a accurate approximation of the underlying eigenvalue of  $T_\mu$ . For  $l \in \mathbb{I}_v^+$ , we have  $\widehat{\lambda}_l \geq (2 - Y_l)\widetilde{\lambda}_l$ , with equality when  $Y_l = 1$ .

**Remark 4.2.** Consider any measure  $\nu \in \mathcal{T}(K)$ ; the approximations  $\widehat{\varphi}_l$ ,  $\vartheta_l/\rho$ ,  $\widehat{\lambda}_l$  and  $\widetilde{\lambda}_l$  (i.e., the normalised approximated eigenfunctions of  $T_\mu$  induced by  $T_\nu$ , the rescaled eigenvalues, and approximated eigenvalues (4.3) and (4.4), respectively) remain unchanged if one replaces  $\nu$  by  $\theta\nu$  for any  $\theta > 0$ .  $\triangleleft$

**5. Analogy with one-class SVM.** Let  $\mathcal{F}$  be the RKHS associated with the squared-kernel  $K^2(\cdot, \cdot)$ . Following [16], problems (3.1) and (3.2) can be interpreted as dual formulations of *squared-kernel one-class distorted SVMs* (or discrepancy-SVMs).

We introduce the function  $g_\omega \in \mathcal{F}$ , defined by  $g_\omega(x) = \sum_{k=1}^N \omega_k K^2(x, x_k)$ , we shall refer to  $g_\omega$  as the *primal distortion term*.

**5.1. One-class SVM related to the regularised problem.** We first describe the SVM related to problem (3.1). For  $g \in \mathcal{F}$ , we consider the convex minimisation problem

$$\begin{aligned} & \underset{g}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{F}}^2 + (g|g_\omega)_{\mathcal{F}} \\ & \text{subject to} && g(x_k) \geq -\alpha d_k \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (5.1)$$

The application  $g \mapsto \|g\|_{\mathcal{F}}^2$  being strictly convex, a solution to problem (5.1) is necessarily unique.

**Lemma 5.1.** *If  $\mathbf{v}_\alpha^*$  is a solution to (3.1) with  $\alpha \geq 0$ , then  $g_\alpha^*(x) = \sum_{k=1}^N [\mathbf{v}_\alpha^* - \boldsymbol{\omega}]_k K^2(x, x_k)$  is the solution to (5.1). For all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ , we have  $g_\alpha^*(x_k) = -\alpha d_k$  (notice that for all  $k$ , we have  $g_\alpha^*(x_k) = [\mathbf{S}(\mathbf{v}_\alpha^* - \boldsymbol{\omega})]_k$ ).*

By introducing the change of variable  $\check{g} = g + g_\omega \in \mathcal{F}$ , problem (5.1) leads to (up to an additive constant)

$$\begin{aligned} & \underset{\check{g}, \gamma}{\text{minimise}} && \frac{1}{2} \|\check{g}\|_{\mathcal{F}}^2 \\ & \text{subject to} && \check{g}(x_k) - g_\omega(x_k) \geq -\alpha d_k \text{ for all } k \in \{1, \dots, N\}, \end{aligned} \quad (5.2)$$

which is an equivalent formulation for (5.1), with solution  $\check{g}_\alpha^*(x) = \sum_{k=1}^N [\mathbf{v}_\alpha^*]_k K^2(x, x_k)$ . Notice that, for all  $x \in \mathcal{X}$ ,  $K^2(x, \cdot) \geq 0$ , so that  $g_\omega$  and  $\check{g}_\alpha^*$  are positive functions on  $\mathcal{X}$ .

**5.2. One-class SVM related to the constrained problem.** We now describe the SVM related to problem (3.2). For  $g \in \mathcal{F}$  and  $\gamma \in \mathbb{R}$ , we introduce the problem

$$\begin{aligned} & \underset{g, \gamma}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{F}}^2 + (g|g_\omega)_{\mathcal{F}} - \gamma \\ & \text{subject to} && g(x_k) \geq \gamma d_k / \varkappa \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (5.3)$$

Again, a solution to problem (5.3) is necessarily unique.

**Lemma 5.2.** *If  $\mathbf{v}_\varkappa^*$  is a solution to (3.2), then  $g_\varkappa^*(x) = \sum_{k=1}^N [\mathbf{v}_\varkappa^* - \boldsymbol{\omega}]_k K^2(x, x_k)$  and  $\gamma_\varkappa^* = (\mathbf{v}_\varkappa^*)^T \mathbf{S}(\mathbf{v}_\varkappa^* - \boldsymbol{\omega})$  is the solution to (5.3). For all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\varkappa^*]_k > 0$ , we have  $g_\varkappa^*(x_k) = \gamma_\varkappa^* d_k / \varkappa$ .*

From Lemma 5.2, we have  $\gamma_\varkappa^* = \|g_\varkappa^*\|_{\mathcal{F}}^2 + (g_\varkappa^*|g_\omega)_{\mathcal{F}}$ . In view of Lemma 3.3, for  $0 < \varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$ , we know that  $\mathbf{v}_\varkappa^*$  is a solution to (3.1) for  $\alpha = -\gamma_\varkappa^* / \varkappa$ , and since  $\alpha \geq 0$ , we therefore have  $\gamma_\varkappa^* \leq 0$ .

**5.3. Soft-margin-type extensions.** Pursuing the analogy with one-class SVMs, we can define soft-margin-type extensions of problems (5.1) and (5.3), i.e, we can consider models where the inequalities appearing in the constraints can, potentially, be violated, the level of violation being penalised. In this section, we only discuss extensions of problem (5.1), but a similar discussion can be done for problem (5.3).

We introduce  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_N)^T \in \mathbb{R}^N$ ; the components of  $\boldsymbol{\xi}$  are referred to as *slack variables*. Instead of considering the constraints  $g(x_k) \geq -\alpha d_k$ , we can consider the relaxed constraints  $g(x_k) \geq -\alpha d_k - \xi_k$ , while penalising the values taken by  $\xi_k$ . The penalisation considered is related to a *loss function*, see for instance [19]. In what follows, we discuss the model obtained for two popular types of loss functions: the (weighted) hinge loss and the (generalised) square loss.

In view of our initial problem (3.1) (regularised squared-kernel-discrepancy minimisation), soft-margin extensions of problem (5.1) appear as tool to further constrain or penalise the measure  $\nu$  (i.e., the vector  $\mathbf{v}$ ) used to approximate the initial measure  $\mu$  (i.e., the vector  $\boldsymbol{\omega}$ ).

*Weighted hinge loss.* Let  $\mathbf{c} \in \mathbb{R}^N$ , with  $\mathbf{c} \geq 0$ ; the soft-margin extension of problem (5.1) corresponding to a weighted hinge loss consists in the problem

$$\begin{aligned} & \underset{g, \boldsymbol{\xi}}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{F}}^2 + (g|g_\omega)_{\mathcal{F}} + \mathbf{c}^T \boldsymbol{\xi} \\ & \text{subject to} && g(x_k) \geq -\alpha d_k - \xi_k, \text{ with } \xi_k \geq 0, \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (5.4)$$

The Lagrange dual of problem (5.4) is given by

$$\underset{\mathbf{v}}{\text{minimise}} D_\alpha(\mathbf{v}) = \frac{1}{2} (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) + \alpha \mathbf{d}^T \mathbf{v} \text{ subject to } 0 \leq \mathbf{v} \leq \mathbf{c}. \quad (5.5)$$

Problem (5.5) only differs from problem (3.1) by the presence of the additional constraints  $\mathbf{v} \leq \mathbf{c}$ , which acts as an upper bound on the values taken by the components of  $\mathbf{v}$  (i.e., an upper bound on the weights of the points in the quadrature related to  $\nu$ ).

*Generalised square loss.* Let  $\boldsymbol{\Sigma} \in \mathbb{R}^{N \times N}$  be a symmetric positive-definite matrix; the soft-margin extension of problem (5.1) corresponding to a generalised square loss consists in the problem

$$\begin{aligned} & \underset{g}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{F}}^2 + (g|g_\omega)_{\mathcal{F}} + \frac{1}{2} \boldsymbol{\xi}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\xi} \\ & \text{subject to} && g(x_k) \geq -\alpha d_k - \xi_k, \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (5.6)$$

The Lagrange dual of problem (5.6) is given by

$$\underset{\mathbf{v}}{\text{minimise}} D_{\alpha, \boldsymbol{\Sigma}}(\mathbf{v}) = \frac{1}{2} (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) + \frac{1}{2} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} + \alpha \mathbf{d}^T \mathbf{v} \text{ subject to } \mathbf{v} \geq 0. \quad (5.7)$$

In comparison to problem (3.1), the term  $\frac{1}{2}\mathbf{v}^T\boldsymbol{\Sigma}\mathbf{v}$  is added to the initial cost  $D_\alpha(\mathbf{v})$ , and tends to “harmonise” the components of the underlying solution  $\mathbf{v}^*$ . In particular, notice that  $D_{\alpha,\boldsymbol{\Sigma}}$  is strongly convex.

**6. Regularisation path.** Considering the framework of Sections 3 and 5 and following the classical results relative to the regularisation paths for Lasso or SVM models (see e.g., [13, 9]), we now discuss the regularisation paths related to problems (3.1) and (3.2). In what follows, we mainly consider problem (3.1) (i.e., the regularised model); results related to problem (3.2) can then be obtained from Lemma 3.3.

**6.1. Generalities.** Let  $\mathbf{v}_\alpha^*$  be a solution to (3.1) for  $\alpha \geq 0$ ; we introduce the index sets

$$J_\alpha = \{k | [\nabla D_\alpha(\mathbf{v}_\alpha^*)]_k = 0\} \text{ and } J_\alpha^c = \{1, \dots, N\} \setminus J_\alpha,$$

so that, by definition,  $[\nabla D_\alpha(\mathbf{v}_\alpha^*)]_{J_\alpha^c} = [\mathbf{S}(\mathbf{v}_\alpha^* - \boldsymbol{\omega}) + \alpha \mathbf{d}]_{J_\alpha^c} > 0$ . From Sections 3 and 5, the index set  $J_\alpha$  is unique, even when the solution to (3.1) is not (i.e., in case of non-uniqueness of the solution,  $J_\alpha$  does not depend on the solution  $\mathbf{v}_\alpha^*$  considered). We shall refer to  $J_\alpha$  as the *sparsity pattern* of the solutions to problem (3.1) for  $\alpha \geq 0$ . We also recall that for all  $k$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ , we have  $k \in J_\alpha$ , see for instance Lemma 5.1.

Knowing  $J_\alpha$ , a solution  $\mathbf{v}_\alpha^*$  is characterised by the conditions

$$[\mathbf{v}_\alpha^*]_{J_\alpha^c} = 0, \text{ and } \mathbf{S}_{J_\alpha, J_\alpha} [\mathbf{v}_\alpha^*]_{J_\alpha} = [\mathbf{S}\boldsymbol{\omega}]_{J_\alpha} - \alpha \mathbf{d}_{J_\alpha} \text{ with } [\mathbf{v}_\alpha^*]_{J_\alpha} \geq 0, \quad (6.1)$$

where  $\mathbf{S}_{J_\alpha, J_\alpha}$  stands for the  $n_\alpha \times n_\alpha$  principal submatrix of  $\mathbf{S}$  corresponding to the index set  $J_\alpha$ , with  $n_\alpha = \text{card}(J_\alpha)$ , and where, for instance,  $\mathbf{d}_{J_\alpha} \in \mathbb{R}^{n_\alpha}$  stands for the vector defined by the components of  $\mathbf{d}$  with index in  $J_\alpha$ . Note that the condition  $[\mathbf{v}_\alpha^*]_{J_\alpha} \geq 0$  is involved only in case of non-uniqueness of the solution.

**Lemma 6.1.** *Let  $\mathbf{v}_{\alpha_1}^*$  and  $\mathbf{v}_{\alpha_2}^*$  be solutions to problem (3.1) with  $\alpha_1$  and  $\alpha_2 \geq 0$ , respectively. Assume that  $J_{\alpha_1} = J_{\alpha_2}$ , then for all  $\theta \in [0, 1]$ ,  $\mathbf{v}_\alpha^* = \theta \mathbf{v}_{\alpha_1}^* + (1 - \theta) \mathbf{v}_{\alpha_2}^*$  is a solution to problem (3.1) with  $\alpha = \theta \alpha_1 + (1 - \theta) \alpha_2$ .*

If, for all  $\alpha \geq 0$ , the solution  $\mathbf{v}_\alpha^*$  is unique (which is for instance the case when  $\mathbf{S}$  is invertible), then in view of Lemma 6.1, the *regularisation map*  $R : \alpha \mapsto \mathbf{v}_\alpha^*$  is a piecewise linear application from  $\mathbb{R}_+$  into  $\mathbb{R}^N$ . In case of uniqueness of the solutions, the regularisation map is therefore piecewise continuous and has right and left limits for all  $\alpha$ , these limits satisfying the optimality conditions for problem (3.1); by uniqueness of the solution, we can in this way show that  $R : \alpha \mapsto \mathbf{v}_\alpha^*$  is continuous. In case of non-uniqueness of the solution, Lemma 6.1 shows that the set of solutions related to a same sparsity pattern  $J$  is convex.

When  $\alpha$  decreases or increases, we refer to a change in the sparsity pattern  $J_\alpha$  as an *event*; the values of  $\alpha$  where an event occurs are sometime referred to as *kinks*. In particular remark that since there can not exist more than  $2^N$  different subsets of  $\{1, \dots, N\}$ , Lemma 6.1 induces that the number  $M_{ev}$  of events related to problem (3.1) necessarily satisfies  $M_{ev} \leq 2^N$ .

In the general case (i.e., even in case of non-uniqueness of the solutions), when  $\alpha$  decreases, we can easily define the value  $\alpha_0$  and the underlying index set  $J_{\alpha_0}$  at which the first event occurs. Indeed, from Lemma 3.2, we know that for  $\alpha \geq \max_k [\mathbf{S}\boldsymbol{\omega}]_k / d_k$ , we have  $\mathbf{v}_\alpha^* = 0$ . Therefore, the first event occurs at  $\alpha_0 = \max_k [\mathbf{S}\boldsymbol{\omega}]_k / d_k$ , and we have  $J_{\alpha_0} = \{k | [\mathbf{S}\boldsymbol{\omega}]_k / d_k = \alpha_0\}$  (for larger values of  $\alpha$ , the underlying sparsity pattern is the empty set). As detailed in the following Section 6.2, if the submatrix  $\mathbf{S}_{J_{\alpha_0}, J_{\alpha_0}}$  is invertible, we can easily compute the value  $\alpha_1$  at which the next event occurs, and also describe the related sparsity pattern  $J_{\alpha_1}$ .

**6.2. Regularisation direction for non-singular submatrix.** We now discuss, in case of uniqueness of the solutions, the computation of the regularisation path for decreasing values of  $\alpha$ , i.e., we assume that the events occur successively at  $\alpha_0 = \max_k [\mathbf{S}\boldsymbol{\omega}]_k / d_k > \alpha_1 > \dots > \alpha_{M_{ev}-1} \geq 0$ , the last event corresponding to the largest value of  $\alpha$  such that  $J_\alpha = \{1, \dots, N\}$ , since for  $\alpha = 0$ , we have  $\mathbf{v}_\alpha^* = \boldsymbol{\omega} > 0$ .

More precisely, considering a kink  $\alpha_p$  with related sparsity pattern  $J_{\alpha_p}$  (for  $p \in \{0, \dots, M_{ev} - 2\}$ ) and assuming that the submatrix  $\mathbf{S}_{J_{\alpha_p}, J_{\alpha_p}}$  is invertible, we describe how to compute the value  $\alpha_{p+1} < \alpha_p$

corresponding to the next event, and how to characterise the related sparsity pattern  $J_{\alpha_{p+1}}$ . We assume that  $n_{\alpha_p} = \text{card}(J_{\alpha_p}) < N$  and, for simplicity, we use the notation  $J = J_{\alpha_p}$ . We recall that, by definition,  $\alpha_p$  is the largest value of  $\alpha$  such that  $J_\alpha = J$ .

From (6.1), we introduce the vector  $\mathbf{v}_\alpha$  such that  $[\mathbf{v}_\alpha]_{J^c} = 0$  and  $[\mathbf{v}_\alpha]_J = \mathbf{S}_{J,J}^{-1}([\mathbf{S}\boldsymbol{\omega}]_J - \boldsymbol{\alpha}\mathbf{d}_J)$ ; the vector  $\mathbf{v}_\alpha$  is sometime referred to as the *regularisation direction*. By definition,  $\alpha_{p+1}$  corresponds to the smallest  $\alpha$  such that  $0 \leq \alpha < \alpha_p$  and

$$[\mathbf{v}_\alpha]_J \geq 0 \text{ and } [\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \boldsymbol{\alpha}\mathbf{d}]_{J^c} \geq 0. \quad (6.2)$$

The set  $J_{\alpha_{p+1}}$  is then obtained by removing from  $J = J_{\alpha_p}$  all the indices  $k \in J$  such that  $[\mathbf{v}_{\alpha_{p+1}}]_k = 0$ , and by adding all the indices  $k \in J^c$  such that  $[\nabla D_{\alpha_{p+1}}(\mathbf{v}_{\alpha_{p+1}})]_k = 0$ ; see Remark 6.1 for more details concerning the computation of  $\alpha_{p+1}$ .

If  $\mathbf{S}_{J_{\alpha_{p+1}}, J_{\alpha_{p+1}}}$  is invertible, we can next compute  $\alpha_{p+2}$  and  $J_{\alpha_{p+2}}$  in exactly the same way, and we may potentially iterate like this until we reach the last event, or, at least, as far as the encountered principal submatrices are invertible.

**Remark 6.1.** Consider (6.2) and introduce the  $(N - n_{\alpha_p}) \times n_{\alpha_p}$  matrix  $\mathbf{M} = \mathbf{S}_{J^c, J} \mathbf{S}_{J, J}^{-1}$ . The condition  $[\mathbf{S}(\mathbf{v}_{\alpha_p} - \boldsymbol{\omega}) + \boldsymbol{\alpha}\mathbf{d}]_{J^c} \geq 0$  yields  $\boldsymbol{\alpha}(\mathbf{M}\mathbf{d}_J - \mathbf{d}_{J^c}) \leq \mathbf{M}[\mathbf{S}\boldsymbol{\omega}]_J - [\mathbf{S}\boldsymbol{\omega}]_{J^c}$ , and we can thus define

$$\alpha_+ = \max_l \left\{ \frac{[\mathbf{M}[\mathbf{S}\boldsymbol{\omega}]_J - [\mathbf{S}\boldsymbol{\omega}]_{J^c}]_l}{[\mathbf{M}\mathbf{d}_J - \mathbf{d}_{J^c}]_l} \mid [\mathbf{M}\mathbf{d}_J - \mathbf{d}_{J^c}]_l < 0 \right\}. \quad (6.3)$$

In the same way, the condition  $[\mathbf{v}_\alpha]_J \geq 0$  gives  $\boldsymbol{\alpha}\mathbf{S}_{J,J}^{-1}\mathbf{d}_J \leq \mathbf{S}_{J,J}^{-1}[\mathbf{S}\boldsymbol{\omega}]_J$  and we obtain

$$\alpha_- = \max_m \left\{ \frac{[\mathbf{S}_{J,J}^{-1}[\mathbf{S}\boldsymbol{\omega}]_J]_m}{[\mathbf{S}_{J,J}^{-1}\mathbf{d}_J]_m} \mid [\mathbf{S}_{J,J}^{-1}\mathbf{d}_J]_m < 0 \right\}. \quad (6.4)$$

The next event then occurs at  $\alpha_{p+1} = \max\{\alpha_+, \alpha_-\}$ . If  $\alpha_{p+1} = \alpha_+$ , the event consists in the entry of new indices in the sparsity pattern, and if  $\alpha_{p+1} = \alpha_-$ , some indices go out of the sparsity pattern.  $\triangleleft$

Once  $\alpha_p$  and  $J = J_{\alpha_p}$  are known, the computation of the next event (i.e., of  $\alpha_{p+1}$  and  $J_{\alpha_{p+1}}$ ) requires the calculation of  $\mathbf{S}_{J,J}^{-1}\mathbf{d}_J$  and  $\mathbf{S}_{J,J}^{-1}[\mathbf{S}\boldsymbol{\omega}]_J$  (i.e., resolution of a linear system). Starting “from scratch” (i.e., without taking into account the computations already performed to obtain the information relative to the kink  $\alpha_p$ ) and using a direct method (by for instance considering the Cholesky decomposition of the symmetric and positive-definite matrix  $\mathbf{S}_{J,J}$ ), the amount of computations required for this operation scale as  $\mathcal{O}(n_{\alpha_p}^3)$ . It is however possible to reduce the computational complexity by considering update formulae (by for instance iteratively updating the Cholesky decomposition of  $\mathbf{S}_{J,J}$ ); in the favorable cases, the computational complexity may thus be reduced to  $\mathcal{O}(n_{\alpha_p}^2)$  (while still considering a direct approach). Notice that an alternative may also consist in considering an indirect iterative approach, like for example conjugate gradient methods. Finally, the complexity of the two matrix-vector products involving the matrix  $\mathbf{S}_{J^c, J}$  scale as  $\mathcal{O}(n_{\alpha_p}(N - n_{\alpha_p}))$ . The computation of the regularisation direction therefore becomes intractable once large values of  $n_{\alpha_p}$  are reached; therefore, when  $N$  is large, the regularisation-path strategy may only be used to explore the range of very sparse approximate measure  $\nu$ .

More generally, the regularisation-path strategy offers an interesting framework to compute solutions to problems (3.1) and (3.2) (or, at least, to obtain very accurate approximations of these solutions). However, the determination of the path is extremely sensitive to numerical errors. For instance, very small value of the gap  $\alpha_p - \alpha_{p+1}$  between two consecutive kinks or the simultaneous entry or exit of indices of the sparsity pattern may lead to numerical precision issues.

See Sections 8 and 9 for applications of the regularisation-path strategy.

**7. Numerical solver for the constrained problem.** In this section, we discuss a strategy to compute approximate solutions to problem (3.2) (i.e., constrained discrepancy minimisation) efficiently and in a numerically tractable way (i.e., without storing the kernel matrix  $\mathbf{S}$ ), for any given value of the parameter  $\kappa \geq 0$ .

**7.1. Vertex-exchange QP solver.** Consider problem (3.2); for  $\varkappa > 0$ , we can define the change of variable  $\tilde{\mathbf{v}} = \mathbf{R}\mathbf{v}$ , with  $\mathbf{R} = \text{diag}(\mathbf{r})$ , and  $\mathbf{r} = (r_1, \dots, r_N)^T = \mathbf{d}/\varkappa$ , i.e.,  $\mathbf{R}$  is a diagonal matrix with  $i$ -th diagonal entry  $r_i$ . In this way, problem (3.2) is turned into (up to an additive constant), for  $\tilde{\mathbf{v}} \in \mathbb{R}^N$ ,

$$\underset{\tilde{\mathbf{v}}}{\text{minimise}} C(\tilde{\mathbf{v}}) = \frac{1}{2}\tilde{\mathbf{v}}^T \mathbf{A}\tilde{\mathbf{v}} - \mathbf{b}^T \tilde{\mathbf{v}} \text{ subject to } \tilde{\mathbf{v}} \geq 0 \text{ and } \mathbf{1}^T \tilde{\mathbf{v}} = 1, \quad (7.1)$$

with  $\mathbf{A} = \mathbf{R}^{-1}\mathbf{S}\mathbf{R}^{-1}$  and  $\mathbf{b} = \mathbf{R}^{-1}\mathbf{S}\boldsymbol{\omega}$ . We refer to (7.1) as the *canonical QP* related to the constrained squared-kernel-discrepancy minimisation (3.2). Since  $\mathbf{A}_{i,j} = K^2(x_i, x_j)/(r_i r_j)$ , any entry of  $\mathbf{A}$  can be easily obtained from only the knowledge of the squared-kernel  $K^2(\cdot, \cdot)$ , the set  $S$  and the vector  $\mathbf{r}$ . We shall therefore not store the matrix  $\mathbf{A}$ , but rather compute on the fly any required entry of  $\mathbf{A}$  (“kernelised solver”). Notice that since  $\boldsymbol{\omega}$  is non-sparse, for large  $N$ , the computation of the *distortion term*  $\mathbf{S}\boldsymbol{\omega}$  is computationally demanding ( $\mathcal{O}(N^2)$  complexity), but it may be parallelised. Once  $\mathbf{b}$  is known, the gradient  $\nabla C(\tilde{\mathbf{v}}) = \mathbf{A}\tilde{\mathbf{v}} - \mathbf{b}$  can be easily obtained for any sparse feasible  $\tilde{\mathbf{v}}$ .

The extreme points of the polytope defined by the constraints in (7.1) are the vectors  $\{\mathbf{e}_i\}_{i=1}^N$ , where  $\mathbf{e}_i \in \mathbb{R}^N$  is the  $i$ -th element of the canonical basis of  $\mathbb{R}^N$  (that is  $[\mathbf{e}_i]_i = 1$ , all the other components being equal to zero). For a feasible  $\tilde{\mathbf{v}}$ , let  $J_{\tilde{\mathbf{v}}} = \{k | \tilde{v}_k > 0\}$  be the index set of all strictly positive components of  $\tilde{\mathbf{v}}$ . An iteration of the vertex-exchange algorithm consists in searching ( $\mathcal{O}(N)$  complexity)

$$i^* = \underset{i}{\text{argmin}}[\nabla C(\tilde{\mathbf{v}})]_i \text{ and } j^* = \underset{j \in J_{\tilde{\mathbf{v}}}}{\text{argmax}}[\nabla C(\tilde{\mathbf{v}})]_j,$$

defining the sparse descent direction  $\boldsymbol{\delta} = \mathbf{e}_{j^*} - \mathbf{e}_{i^*}$  (i.e., weight is transferred from the  $j^*$ -th to the  $i^*$ -th component of  $\tilde{\mathbf{v}}$ ); in case of non-uniqueness of the extrema, an index is simply selected at random among the ones satisfying the condition. The step size is then classically obtained by line search, the optimal step size  $\sigma$  being given by  $\sigma = \min\{\tilde{v}_{j^*}, -(\boldsymbol{\delta}^T \nabla C(\tilde{\mathbf{v}}))/(\boldsymbol{\delta}^T \mathbf{A}\boldsymbol{\delta})\}$ . In particular, since the descent direction  $\boldsymbol{\delta}$  is sparse, the computation of the optimal step size is numerically inexpensive, and the same holds for the gradient update. Indeed, we have  $\nabla C(\tilde{\mathbf{v}} + \sigma\boldsymbol{\delta}) = \nabla C(\tilde{\mathbf{v}}) + \sigma\mathbf{A}\boldsymbol{\delta}$ , so that the gradient update only involves two columns of  $\mathbf{A}$ . The kernelised vertex-exchange strategy thus appears as an interesting candidate to scale up to relatively large  $N$  and may be used as a complement of the regularisation-path strategy described in Section 6.

Denoting by  $\tilde{\mathbf{v}}^*$  a solution to (7.1), the convergence of the vertex-exchange algorithm can be easily verified (see, e.g., [7]) by simply remarking that since  $\tilde{\mathbf{v}} \geq 0$  and  $\mathbf{1}^T \tilde{\mathbf{v}} = 1$ , by definition of  $j^*$ , we have  $\tilde{\mathbf{v}}^T \nabla C(\tilde{\mathbf{v}}) \leq \mathbf{e}_{j^*}^T \nabla C(\tilde{\mathbf{v}})$ , so that

$$C(\tilde{\mathbf{v}}) - C(\tilde{\mathbf{v}}^*) \leq -(\mathbf{e}_{i^*} - \tilde{\mathbf{v}})^T \nabla C(\tilde{\mathbf{v}}) \leq -(\mathbf{e}_{i^*} - \mathbf{e}_{j^*})^T \nabla C(\tilde{\mathbf{v}}),$$

and these inequalities can also be used to check distance from optimality. In Sections 8 and 9, the accuracy of an approximate solution  $\tilde{\mathbf{v}}$  is indicated by  $\varepsilon = (\tilde{\mathbf{v}} - \mathbf{e}_{i^*})^T \nabla C(\tilde{\mathbf{v}})$  (Frank-Wolfe error bound).

**7.2. Enhancing sparsity through components merging.** The canonical QP formulation introduced in Section 7.1 offers a convenient framework to attempt to enhance the sparsity of the approximation while trying to keep the discrepancy of the related measure  $\nu$  as low as possible. Let  $\tilde{\mathbf{v}} \geq 0$  be such that  $\mathbf{1}^T \tilde{\mathbf{v}} = 1$ . In practice,  $\tilde{\mathbf{v}}$  will be an exact or approximate solution to problem (7.1), or any vector related to an interesting low-discrepancy configuration  $\mathbf{v}$  through the change of variable  $\tilde{\mathbf{v}} = \mathbf{R}\mathbf{v}$ , with  $\mathbf{R} = \text{diag}(\mathbf{r})$  and  $\mathbf{r} = \mathbf{d}/(\mathbf{d}^T \mathbf{v})$ , see Section 7.1. We assume that  $\tilde{\mathbf{v}}$  has  $n = n_0$  strictly positive components and we introduce  $I = \{i | \tilde{v}_i > 0\}$ . In practice (see Section 8), we observe that it appears possible, to a certain extent, to merge together some components of  $\tilde{\mathbf{v}}$  while inducing a negligible increase of the cost  $C(\cdot)$ . In what follows, we discuss two simple greedy heuristics based on the sequential merging of pairs of components of  $\tilde{\mathbf{v}}$ .

We assume that  $n > 1$ . For an ordered pair  $\{i, j\}$ , with  $i$  and  $j \in I$  and  $i \neq j$ , we define  $\tilde{\mathbf{v}}_{\{i,j\}} = \tilde{\mathbf{v}} + \tilde{v}_j(\mathbf{e}_i - \mathbf{e}_j)$ , i.e.,  $\tilde{\mathbf{v}}_{\{i,j\}}$  has  $n - 1$  strictly positive components, the  $j$ -th component of  $\tilde{\mathbf{v}}$  being absorbed by the  $i$ -th; we refer to this operation as the  $\{i, j\}$ -merging of  $\tilde{\mathbf{v}}$ . We have

$$C(\tilde{\mathbf{v}}_{\{i,j\}}) = C(\tilde{\mathbf{v}}) + \frac{1}{2}\tilde{v}_j^2(\mathbf{e}_i - \mathbf{e}_j)^T \mathbf{A}(\mathbf{e}_i - \mathbf{e}_j) + \tilde{v}_j(\mathbf{e}_i - \mathbf{e}_j)^T \nabla C(\tilde{\mathbf{v}}).$$

Thus, knowing  $\nabla C(\tilde{\mathbf{v}})$ , the computation  $C(\tilde{\mathbf{v}}_{\{i,j\}})$  is numerically inexpensive (since only four entries of the matrix  $\mathbf{A}$  and two entries of  $\nabla C(\tilde{\mathbf{v}})$  are involved).

We can then search for the merging associated with the smallest value of  $C(\tilde{\mathbf{v}}_{\{i,j\}})$ , with  $i$  and  $j \in I$ , and  $i \neq j$ . Depending on  $n_0$  and on the computational power at disposal, we may either consider

- *strong-pairwise-merging*: search for the best ordered pair  $\{i^*, j^*\} = \operatorname{argmin}_{i \neq j} = (\tilde{\mathbf{v}}_{\{i,j\}})$ , the amount of computations involved scaling as  $\mathcal{O}(n^2)$ ; or
- *weak-pairwise-merging*: fix  $j^* = \operatorname{argmin}_{j \in I} \tilde{v}_j$ , and search for  $i^* = \operatorname{argmin}_{i \neq j^*} = (\tilde{\mathbf{v}}_{\{i,j^*\}})$ , the amount of computations involved scaling as  $\mathcal{O}(n)$ .

We thus obtain a “best” pairwise merging  $\{i^*, j^*\}$  for  $\tilde{\mathbf{v}}$ . We next update all the involved objects, i.e.,  $\tilde{\mathbf{v}} \leftarrow \tilde{\mathbf{v}}_{\{i^*, j^*\}}$ ,  $I \leftarrow I \setminus \{j^*\}$ ,  $n \leftarrow n - 1$  and  $\nabla C(\tilde{\mathbf{v}}) \leftarrow \nabla C(\tilde{\mathbf{v}}_{\{i^*, j^*\}})$ , and we may potentially iterate like this until  $n = 1$  (i.e., after  $n_0 - 1$  iterations), or at least, until we have reached a satisfactory sparsity-discrepancy trade-off.

We thus obtain a sequence of merged-vectors  $\{\tilde{\mathbf{v}}_{[0]}, \tilde{\mathbf{v}}_{[1]}, \dots, \tilde{\mathbf{v}}_{[n_0-1]}\}$ , where  $\tilde{\mathbf{v}}_0$  is our initial vector,  $\tilde{\mathbf{v}}_{[1]}$  results from the merging of two components of  $\tilde{\mathbf{v}}_{[0]}$ , etc.; by construction,  $\tilde{\mathbf{v}}_{[k]} \geq 0$  and  $\mathbb{1}^T \tilde{\mathbf{v}}_{[k]} = 1$  for all  $k$ , and  $\tilde{\mathbf{v}}_{[k]}$  has  $n_0 - k$  strictly positive components. Finally, instead of considering the Nyström approximation induced by  $\mathbf{v} = \mathbf{R}^{-1} \tilde{\mathbf{v}}_0$ , we may consider a sparser vector  $\mathbf{v}_{[k]} = \mathbf{R}^{-1} \tilde{\mathbf{v}}_{[k]}$  (notice that  $\mathbf{v}_{[k]}$  and  $\tilde{\mathbf{v}}_{[k]}$  have the same number of strictly positive components); see Sections 8 and 9 for illustrations.

**8. Two-dimensional example.** We assume that  $S = \{x_k\}_{k=1}^N$  consists of the  $N = 2016$  first points of a uniform Halton sequence on  $[-1, 1]^2$  (see Figure 8.2), and we set  $\omega_k = 1/N$  for all  $k$ ; in particular, the measure  $\mu = \sum_k \omega_k \delta_{x_k}$  thus appears as a quadrature approximation of the uniform probability measure on  $[-1, 1]^2$ . We consider the Gaussian kernel  $K(x, y) = \exp(-\ell \|x - y\|^2)$ , where  $\|x - y\|$  is the Euclidean norm on  $\mathbb{R}^2$ , and we set  $\ell = 1/0.16$ . An overview of the spectrum of the operator  $T_\mu$  is given in Figure 8.1. In all Section 8, we consider  $\mathbf{d} = \operatorname{diag}(\mathbf{K}) = \mathbb{1}$ , i.e., we penalise the trace of  $T_\nu$ ; we have  $\mathbf{d}^T \boldsymbol{\omega} = \operatorname{trace}(T_\mu) = 1$ .

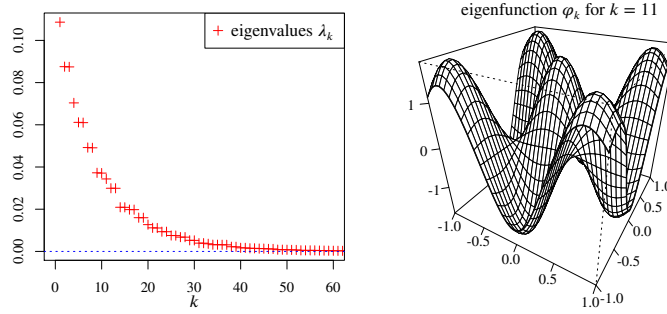


Fig. 8.1: For the two-dimensional example, eigenvalues  $\lambda_k$  of the integral operator  $T_\mu$  (sorted in decreasing order, only the 62 largest eigenvalues are presented), and graph, on  $[-1, 1]^2$  of the canonically extended eigenfunction  $\varphi_k$  for  $k = 11$ .

**8.1. First experiment.** Figure 8.2 shows the (approximate) solution  $\mathbf{v}^*$  to problem (3.2) with  $x = 0.81$ , or equivalently, to problem (3.1) with  $\alpha \approx 8.354215 \times 10^{-3}$ . The vector  $\mathbf{v}^*$  has 160 strictly positive components, and the support of the related measure  $\nu^*$  inherits an interesting “four-concentric-square” structure. We have  $D(\mathbf{v}^*) = 7.631890 \times 10^{-4}$  (for comparison, notice that  $D(\mathbf{0}) = \frac{1}{2} \boldsymbol{\omega}^T \mathbf{S} \boldsymbol{\omega} = 2.661452 \times 10^{-2}$  and  $D(\mathbf{e}_1) = 4.760566 \times 10^{-1}$ , with  $\mathbf{e}_1$  the first element of the canonical basis of  $\mathbb{R}^N$ ).

The presented solution has been obtained using the regularisation-path strategy described in Section 6 (see Section 8.2 for more details). Considering the regularisation path for problem (3.1)

with decreasing values of  $\alpha$ , the underlying value of  $\alpha \approx 8.354215 \times 10^{-3}$  satisfies

$$\alpha_{p+1} = 8.352970 \times 10^{-3} \leq \alpha \leq \alpha_p = 8.355244 \times 10^{-3}, \text{ with } p = 4047.$$

Correspondingly, considering the regularisation path for problem (3.2) with increasing values of  $\varkappa$ , the underlying value  $\varkappa = 0.81$  satisfies

$$\varkappa_p = 0.8099788 \leq \varkappa \leq \varkappa_{p+1} = 0.8100256, \text{ with } p = 4047.$$

In the framework of Section 7.1, the presented solution is related to a Frank-Wolfe error bound  $\epsilon = 4.510281 \times 10^{-17}$ .

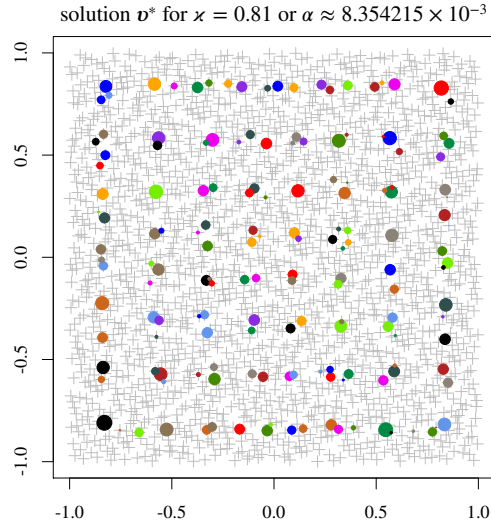


Fig. 8.2: Graphical representation (two-dimensional example) of the solution  $\mathbf{v}^*$  to problem (3.2) with  $\varkappa = 0.81$ , or equivalently, to problem (3.1) with  $\alpha \approx 8.354215 \times 10^{-3}$  (for  $\mathbf{d} = \text{diag}(\mathbf{K})$  and  $\boldsymbol{\omega} = \mathbf{1}/N$ ). The grey crosses represent the points in  $S$  and the filled dots are the strictly positive components of  $\mathbf{v}^*$  (surface being proportional to  $v_k^*$ ).

The efficiency of the Nyström approximation induced by the solution  $\mathbf{v}^*$  presented in Figure 8.2 is illustrated in Figure 8.3. In view of the values of  $Y_k$  (see Section 4), we observe that the 21 main eigendirections of the operator  $T_{\mathbf{v}^*}$  leads to remarkably accurate approximations of the eigenpairs of  $T_\mu$  related to the largest eigenvalues  $\lambda_k$  (for  $k \in \{1, \dots, 21\}$ , we indeed have  $0.9876064 \leq Y_k \leq 0.9999785$ ). The quality of the approximate eigenpairs  $(\hat{\lambda}_k, \hat{\varphi}_k)$  starts to decrease for  $k \geq 22$ . For  $k \in \{22, \dots, 31\}$  the approximate eigenpairs are still relatively accurate; the approximation efficiency decreases significantly for  $k \in \{32, \dots, 43\}$  and becomes very poor for  $k \geq 44$ .

**8.2. Regularisation path.** We now investigate the impact of the parameters  $\alpha$  and  $\varkappa$  related to problems (3.1) and (3.2) respectively. We compute the 12 786 first events of the regularisation-path related to problem (3.1) with decreasing values of  $\alpha$  (see Section 6), i.e., until we reach a precision issue; in particular, we have  $\alpha_0 = 6.310163 \times 10^{-2}$  and  $\alpha_{12785} = 1.514626 \times 10^{-5}$ . Correspondingly, for the regularisation path related to problem (3.1) with increasing  $\varkappa$ , we have  $\varkappa_0 = 0$  and  $\varkappa_{12785} = 0.9995426$  (we recall that  $\mathbf{d}^T \boldsymbol{\omega} = 1$ ).

Figure 8.4 shows that the number of strictly positive components of the solution  $\mathbf{v}_\varkappa^*$  to problem (3.1) tends to increase when  $\varkappa$  increases. As expected from Lemma 3.2-(vi), the functions  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  is decreasing; in the same way, when  $\varkappa$  increases, the corresponding value of the regularisation parameter  $\alpha$  decreases (see Lemmas 3.2 and 3.3).

For 41 values of  $\varkappa$  evenly spread between 0 and  $\varkappa_{12785}$ , Figure 8.5 shows the evolution of the Upsilon-test (4.2) related to the Nyström approximation induced by the various solutions  $\mathbf{v}_\varkappa^*$ . As

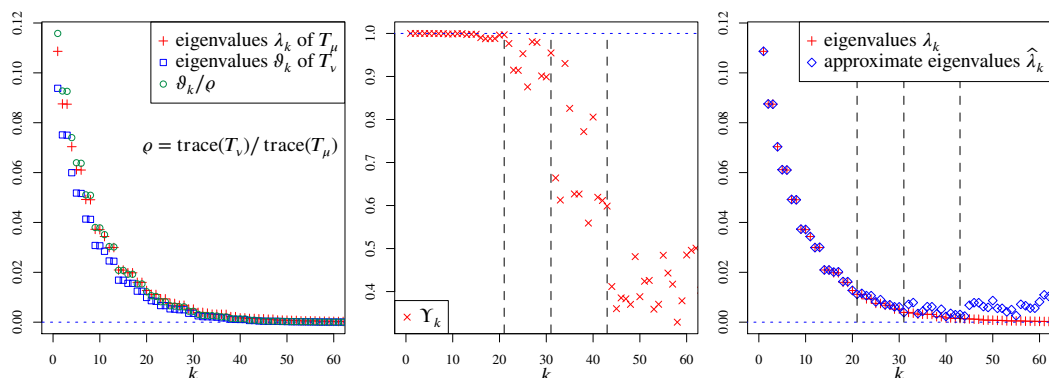


Fig. 8.3: Eigenpairs approximation induced by the solution  $\mathbf{v}^*$  presented in Figure 8.2: graphical representations of the 62 largest eigenvalues  $\lambda_k$ ,  $\vartheta_k$  and rescaled eigenvalue  $\vartheta_k/\rho$ , sorted in decreasing order (left); Upsilon-test (4.2) for the 62 main eigendirections of  $T_{\mathbf{v}^*}$  (middle); and 62 first approximate eigenvalues  $\hat{\lambda}_k$  induced by  $\mathbf{v}^*$  (right).

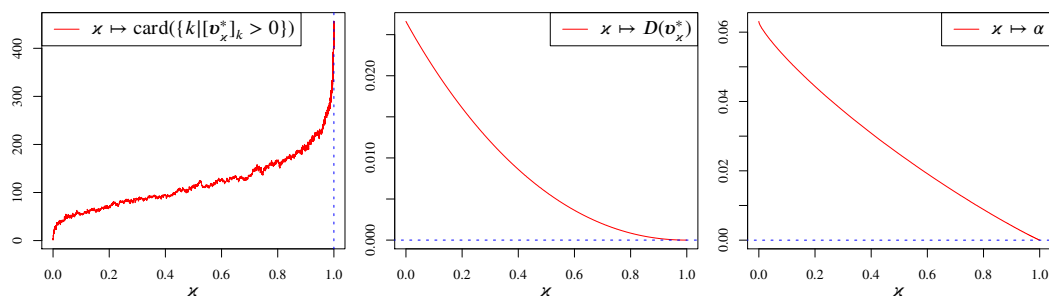


Fig. 8.4: Graphical representation of the 12786 first events of the regularisation path related to problem (3.1) for increasing  $x$  (two-dimensional example): number of strictly positive components in  $\mathbf{v}_x^*$  as function of  $x$  (left); graph of  $x \mapsto D(\mathbf{v}_x^*)$  (middle), and relation between  $x$  and the regularisation parameter  $\alpha$  of problem (3.1) (see Lemma 3.3).

expected, the number of accurately approximate eigendirections increases with  $x$ . Remarkably, for each of the considered values of  $x$ , the number of eigendirections approximated with a high accuracy appears to be in close relation with the decay of the spectrum of  $T_\mu$ ; we recall that, since  $\mathbf{d} = \text{diag}(\mathbf{K})$ , we have  $\text{trace}(T_{\mathbf{v}_x^*}) = x$ .

**8.3. Components merging.** We now perform the strong-pairwise-merging (see Section 7.2) of the solution  $\mathbf{v}^*$  presented in Figure 8.2 (i.e., problem (3.2) with  $x = 0.81$ ). As illustrated in Figure 8.6, for the first merging iterations,  $D(\mathbf{v}_{[k]})$  stays very close to  $D(\mathbf{v}^*) = 7.631890 \times 10^{-4}$ . After 90 iterations, we have  $D(\mathbf{v}_{[90]}) - D(\mathbf{v}^*) = 3.494809 \times 10^{-5}$  (i.e., increasing of 4.58%), and  $\mathbf{v}_{[90]}$  is supported by 70 points (instead of 160 for  $\mathbf{v}^*$ ); a graphical representation of  $\mathbf{v}_{[90]}$  is given in the left-hand part of the figure. The efficiency of the Nyström approximation induced by  $\mathbf{v}_{[90]}$  is presented in the right-hand part of Figure 8.6 (Upsilon-test). We observe that although being slightly less accurate than the approximate eigendecomposition induced by  $\mathbf{v}^*$ , the approximation induced by  $\mathbf{v}_{[90]}$  remains very satisfactory while being related to a vector more than two times sparser.

**8.4. Comparison with random sampling.** For comparison purpose, we compute the approximate eigendecompositions induced by random uniform samples (without repetition) of size  $n_{rand} = 300, 600, 900$  and  $1200$  (i.e., we randomly select  $n_{rand}$  distinct points among the  $N = 2016$  points in  $\mathcal{S}$ , and we consider the uniform probability measure supported by the points selected); for each sample



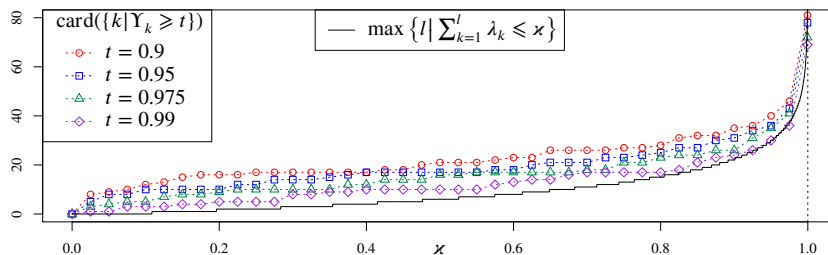


Fig. 8.5: Evolution of the accuracy of the approximate eigendecomposition of  $T_\mu$  induced by  $\mathbf{v}_x^*$  for 41 values of  $x$  between  $x_0 = 0$  and  $x_{12785} = 0.9995426$ ; the accuracy of the approximate eigendirections is measured through the Upsilon-test, and for illustration purpose, the map  $x \mapsto \max \{l \mid \sum_{k=1}^l \lambda_k \leq x\}$  is also presented (two-dimensional example).

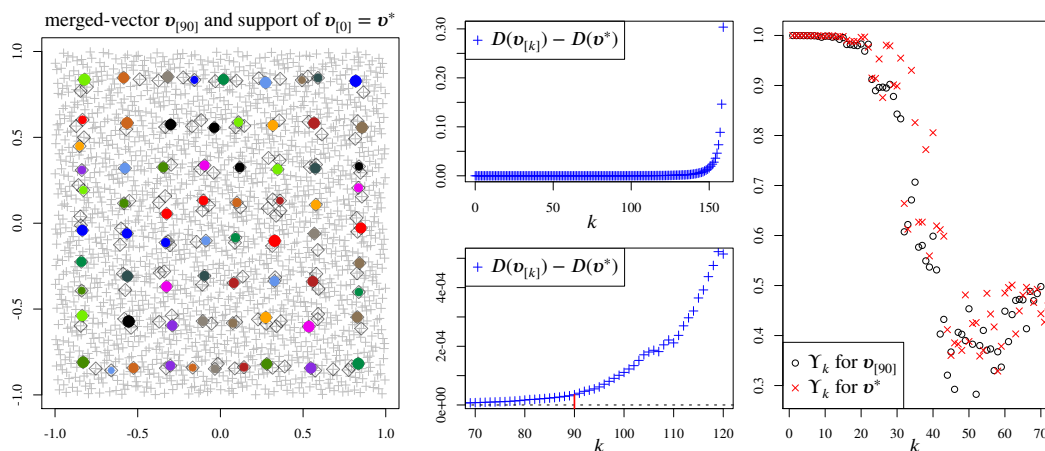


Fig. 8.6: Graphical representation of the merged-solution  $\mathbf{v}_{[90]}$  (two-dimensional example) obtained after 90 iterations of the strong-pairwise-merging strategy applied to the solution  $\mathbf{v}^*$  presented in Figure 8.2; the grey diamonds indicate the support of  $\mathbf{v}^*$  (left). Increase of the cost  $D(\cdot)$  induced by each merging iteration, for the whole 159 iterations (top-middle), and zoom around the 90-th iteration (bottom-middle). Representation of the Upsilon-test obtained from the merged-vector  $\mathbf{v}_{[90]}$  and comparison with the Upsilon-test related to the solution  $\mathbf{v}^*$  (right).

size, we perform 100 repetitions. Figure 8.7 illustrates the accuracy of the resulting approximate eigendirections (measured through the Upsilon-test). As we could expect, the accuracy of the approximations increases with the size of the sample; however, in terms of trade-off between sparsity and accuracy, the results are far behind the one obtained using squared-kernel-discrepancy minimisation.

**9. Application to medium/large-scale problems.** The main motivation behind this section is to illustrate the ability of the proposed framework to tackle relatively large-scale problems. The datasets have been downloaded from the UCI Machine Learning Repository, see [12]. All the computations have been performed on a 2015 desktop endowed with an Intel Core i7-4790 processor with 16 Gb of RAM; the various methods have been entirely encoded in C (the code is available upon request).

**9.1. MiniBooNE dataset.** We consider the standardised entries of the MiniBooNE dataset, without labels;  $\mathcal{S}$  thus consists of  $N = 129\,596$  points in  $\mathbb{R}^{50}$ . We use a Gaussian kernel (same expression as in Section 8) with  $\ell = 0.02$ , and we set  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{d} = \text{diag}(\mathbf{K}) = \mathbf{1}$  (in particular,  $\mathbf{d}^T \boldsymbol{\omega} = 1$ ). Notice that  $\ell = 0.02$  belongs to the range of “good parameters” for the SVM binary classification of this dataset.

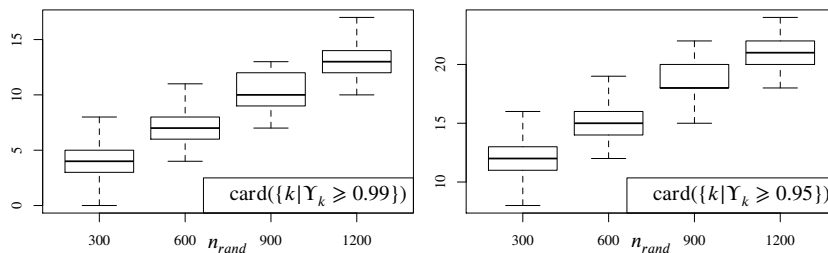


Fig. 8.7: For the two-dimensional example, accuracy of the approximate eigendecompositions induced by random samples of size  $n_{rand}$ ; for each values of  $n_{rand}$ , boxplot (minimum, first quartile, median, third quartile and maximum), over 100 repetitions, of the number of approximate eigendirections such that  $Y_k \geq 0.99$  (left) and  $Y_k \geq 0.95$  (right).

We compute the 3 000 first events of the regularisation path related to problems (3.1) and (3.2). We have  $\alpha_0 = 0.2188961$  and  $\alpha_{2999} = 3.546703 \times 10^{-3}$ , and correspondingly  $\varkappa_0 = 0$  and  $\varkappa_{2999} = 0.655808$ . A graphical representation of the obtained results is proposed in Figure 9.1. We can observe that for  $\varkappa \geq 0.5$ , the number of strictly positive components of the solution  $\mathbf{v}_\varkappa^*$  increases with a significant rate; as a consequence, the computation of the regularisation path quickly becomes numerically intractable (notice that the calculation of the 3 000 first events of the regularisation path has required around 4 hours on our aforementioned 2015 desktop).

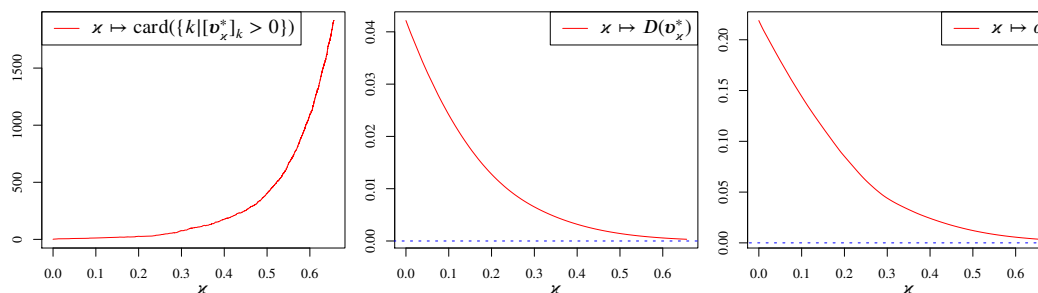


Fig. 9.1: For the MiniBooNE dataset, graphical representation of the 3 000 first events of the regularisation path related to problem (3.1) for increasing  $\varkappa$ : number of strictly positive components in  $\mathbf{v}_\varkappa^*$  as function of  $\varkappa$  (left); graph of  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  (middle), and relation between  $\varkappa$  and the regularisation parameter  $\alpha$  of problem (3.1) (see Lemma 3.3).

We then compute the approximate eigendecompositions induced by the solutions to problem (3.1) for  $\varkappa = 0.3$  and  $\varkappa = 0.655$  obtained with the regularisation-path strategy (or equivalently, to problem (3.1) with  $\alpha \approx 4.400276 \times 10^{-2}$  and  $\alpha \approx 3.571413 \times 10^{-3}$ , respectively). The underlying solutions have 76 and 1 902 strictly positive components, respectively. The efficiency of the induced approximate eigendecompositions is illustrated in Figure 9.2. For  $\varkappa = 0.3$ , we obtain an accurate approximation of the three main eigenpairs of  $T_\mu$  while considering only 76 points (we recall that  $N = 129\,596$ ); the approximation of the other eigendirections is relatively poor. For  $\varkappa = 0.655$ , the eight main eigendirections of  $T_\mu$  are approximate with high accuracy, and the approximation becomes poor for  $k \geq 30$ . We in particular observe a relatively fast decay of the spectrum.

To explore the type of solutions obtained for larger values of  $\varkappa$ , we consider the vertex-exchange strategy described in Section 7.1. We compute an approximate solution for  $\varkappa = 0.8$ ; the vertex-exchange algorithm is initialised at  $\tilde{\mathbf{v}} = \mathbf{e}_1$  and after 300 000 iterations, we obtain a Frank-Wolfe error bound of  $\epsilon = 1.692408 \times 10^{-8}$ . The obtained approximate solution  $\hat{\mathbf{v}}^*$  to problem (3.1) verifies  $D(\hat{\mathbf{v}}^*) = 4.934072 \times 10^{-5}$  and has 9544 strictly positive components. The efficiency of the

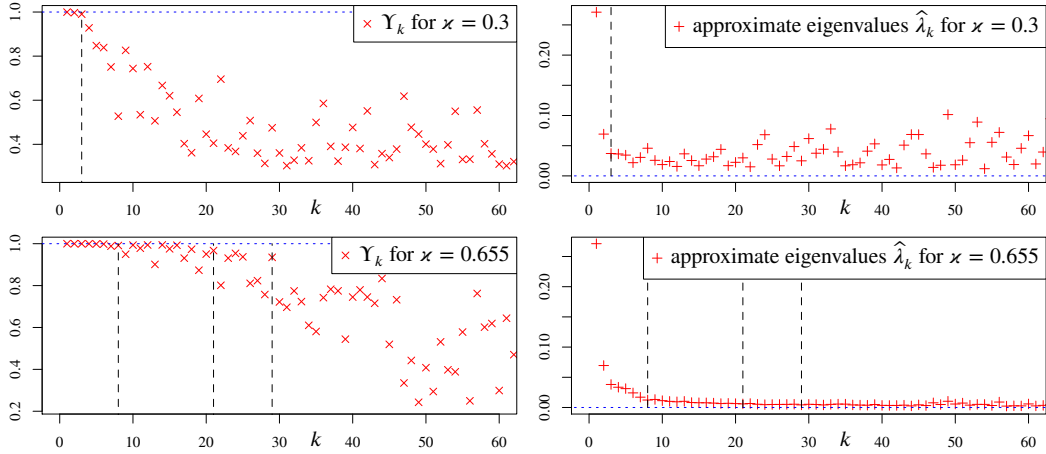


Fig. 9.2: For the MiniBooNE dataset, eigenpairs approximation induced by the solution to problem (3.1) with  $\kappa = 0.3$  (top) and  $\kappa = 0.655$  (bottom): Upsilon-test (4.2) for the 62 main eigendirections of the approximate operator  $T_{v^*}$  (left); and 62 first approximate eigenvalues  $\hat{\lambda}_k$  induced by  $v^*$  (right).

approximate eigendecomposition induced by  $\hat{d}^*$  is illustrated in Figure 9.3; we observe that the 21 main eigendirections of  $T_\mu$  are approximated with high accuracy, and that the 43 main eigendirections are relatively well approximated.

In order to enhance the sparsity of the sampling, we perform a weak-pairwise merging of the solution  $\hat{d}^*$  (see Section 7.2); notice that performing a strong-pairwise merging is in this case numerically prohibitive. After 5044 iterations, the merged solution  $v_{[5044]}$  is supported by 4500 points and  $D(v_{[5044]}) = D(\hat{d}^*) + 1.061787 \times 10^{-6}$  (i.e., increasing of 2.15%). We next compute the approximate eigendecomposition induced by  $v_{[5044]}$ , the result is presented in Figure 9.3. We observe that, in terms of accuracy, the approximation of the main eigendirections of  $T_\mu$  induced by  $v_{[5044]}$  is equivalent to the approximation induced by  $\hat{d}^*$  (and, remarkably, is even substantially better), while being more than two times sparser.

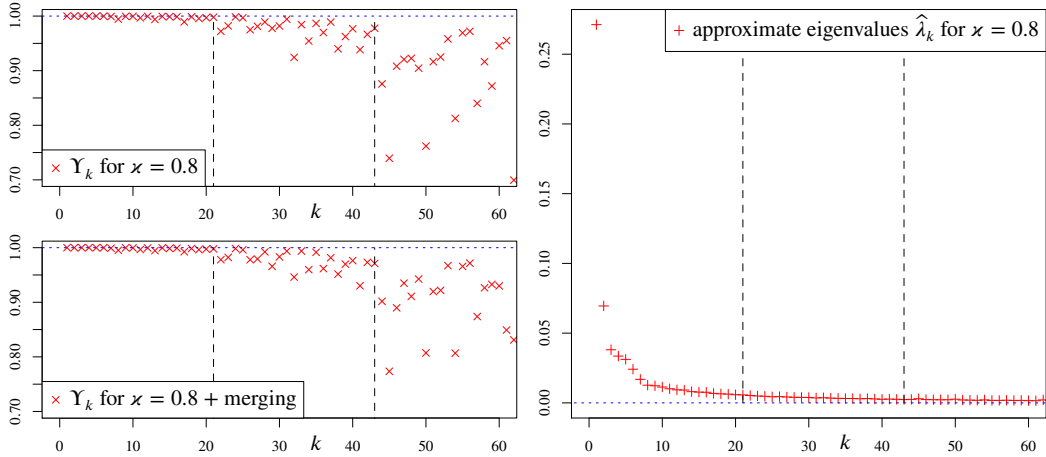


Fig. 9.3: For the MiniBooNE dataset, accuracy of the approximate eigendecompositions (Upsilon-test) induced by the solution  $\hat{d}^*$  to problem (3.1) with  $\kappa = 0.8$  obtained from the vertex-exchange algorithm (top-left) and from the merged solution  $v_{[5044]}$  (bottom-left); graphical representation of the 62 first approximate eigenvalues  $\hat{\lambda}_k$  induced by  $\hat{d}^*$  (right).

To conclude this example, for the various approximations considered in Section 9.1, Table 9.1 presents the values taken by  $\hat{\tau}_m = \sum_{k=1}^m \hat{\lambda}_k$ , where, for each approximations,  $m$  is the number of accurately approximated eigendirections highlighted in Figures 9.2 and 9.3. We can in particular observe that the values of  $\hat{\tau}_m$  obtained from the various approximations are in agreement.

Table 9.1: For the various approximate eigendecompositions considered in Section 9.1 (MiniBooNE dataset), value of  $\sum_{k=1}^m \hat{\lambda}_k$  for specific values of  $m$ .

	$m = 3$	$m = 8$	$m = 21$	$m = 43$
$\varkappa = 0.3$	0.3773318			
$\varkappa = 0.655$	0.3787624	0.4971454	0.6085483	
$\hat{\mathbf{d}}^*$ ( $\varkappa = 0.8$ )	0.3787608	0.4972059	0.6040054	0.6831240
$\mathbf{v}_{[5044]}$ ( $\varkappa = 0.8$ )	0.3787608	0.4972103	0.6039483	0.6829062

**9.2. Test subsample of the SUSY dataset.** We consider the standardised entries of the test subsample of the SUSY dataset (without labels), so that  $\mathcal{S}$  consists of  $N = 500\,000$  points in  $\mathbb{R}^{18}$ . We still use a Gaussian kernel (same expression as in Section 8) with  $\ell = 0.4$ , and we set  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{d} = \text{diag}(\mathbf{K}) = \mathbf{1}$  (so that  $\mathbf{d}^T \boldsymbol{\omega} = 1$ ). The computation of the distortion term  $\mathbf{S}\boldsymbol{\omega}$  takes 5 665.6 seconds.

We compute an approximate solution for the constrained problem (3.2) with  $\varkappa = 0.3$ ; we perform four consecutive batches of 50 000 iterations each, the solver being initialised at  $\tilde{\mathbf{v}} = \mathbf{e}_1$ . After 200 000 iterations (i.e., at the end of the 4-th batch), the obtained approximate solution  $\hat{\mathbf{d}}^*$  verifies  $D(\hat{\mathbf{d}}^*) = 3.931629 \times 10^{-5}$  and has  $n = 20\,664$  strictly positive components. Execution time, evolutions of the Frank-Wolfe error bound  $\epsilon$  and of the sparsity of the approximate solution are reported in Table 9.2. We observe that a batch of 50 000 iterations of the vertex-exchange algorithm takes around 19 minutes, and that the approximate solution obtained at the end of the first batch is already relatively accurate.

Table 9.2: For the test subsample of the SUSY dataset, information relative to approximate solution to problem (3.2) with  $\varkappa = 0.3$  returned by the vertex-exchange algorithm for four consecutive batches of 50 000 iterations, the solver being initialised at  $\tilde{\mathbf{v}} = \mathbf{e}_1$ : execution time, total number of iterations, Frank-Wolfe error bound  $\epsilon$  and number  $n$  of strictly positive component of the approximate solution.

	batch 1	batch 2	batch 3	batch 4
time (in sec.)	1 148.7	1 158.3	1 158.5	1 159.1
total nb. of it.	50 000	100 000	150 000	200 000
$\epsilon$	$3.1413 \times 10^{-7}$	$6.5477 \times 10^{-8}$	$2.7049 \times 10^{-8}$	$7.0928 \times 10^{-9}$
$n$	19 721	20 619	20 693	20 674

In order to enhance the sparsity of the sampling, we perform a weak-pairwise merging of the approximate solution  $\hat{\mathbf{d}}^*$ ; the computation of 20 673 merging iterations takes 78.86 seconds. The merged solution  $\mathbf{v}_{[13674]}$  is supported by 7 000 points and  $D(\mathbf{v}_{[13674]}) = D(\hat{\mathbf{d}}^*) + 5.271960 \times 10^{-7}$  (i.e., increasing of only 1.34%). We next study the approximate eigendecomposition induced by  $\mathbf{v}_{[13674]}$ . Computing the 300 first normalised approximate eigenvectors of  $\mathbf{KW}$  induced by  $\mathbf{v}_{[13674]}$  (see Section 4) takes 3 313.6 seconds (time for canonical extension and rescaling); then, performing the Upsilon-test for this 300 approximate eigendirections takes 191 622.3 seconds (i.e., around 53 hours). The results are presented in Figure 9.4. As already observed, the accuracy of the approximate eigendirections  $\hat{\varphi}_k$  decreases when  $k$  increases (we recall that the eigenvalues of the approximate operator are stored in descending order). The obtained approximate eigendirections are remarkably

accurate: for instance, we have  $Y_k \geq 0.99$  for 64 eigendirections,  $Y_k \geq 0.95$  for 224 eigendirections, and  $Y_k \geq 0.9$  for 299 eigendirections (among 300), and  $\min_{k=1}^{300} Y_k = 0.8826657$ . The approximate eigenvalues  $\hat{\lambda}_k$  and the related eigenvalues  $\vartheta_k$  of the operator  $T_\nu$  defined from  $\mathbf{v}_{[13674]}$  are also presented.

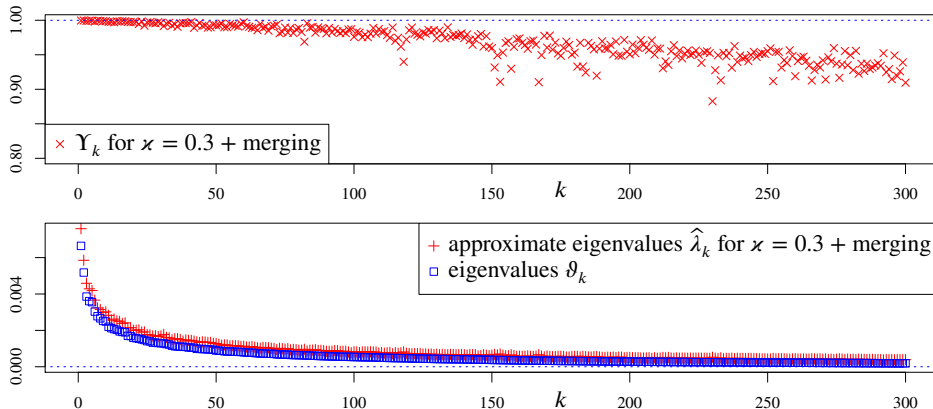


Fig. 9.4: For the test subsample of the SUSY dataset, accuracy of the 300 main approximate eigendirections (Upsilon-test) induced by the merged vector  $\mathbf{v}_{[13674]}$  obtained from the approximate solution  $\hat{\mathbf{d}}^*$  to problem (3.2) with  $x = 0.3$  (top); graphical representation of the underlying approximate eigenvalues  $\hat{\lambda}_k$  and of the corresponding eigenvalues  $\vartheta_k$  of the operator related to  $\mathbf{v}_{[13674]}$  (bottom).

**10. Conclusion.** We proposed a QP-based strategy for computing sparse Nyström approximations of integral operators defined from a symmetric positive-semidefinite kernel and a discrete measure with finite support. Our criterion is based on the notion of squared-kernel discrepancy (i.e., on the squared Hilbert-Schmidt norm for operators defined on the RKHS related to the underlying kernel). For a fixed kernel and starting from a given discrete measure  $\mu$ , we compute an approximate measure  $\nu$  which is as close as possible to  $\mu$  in terms of squared-kernel discrepancy while being supported by a smaller number of points, the support of  $\nu$  being included in the support of  $\mu$  (quadrature-sparsification). From a methodological point of view, the obtained approximation scheme is deterministic, numerically scalable (i.e., large-scale problem can be tackled) and enjoys an optimality property.

In this work, the sparsity of  $\nu$  is enforced through a weighed- $\ell_1$ -type penalisation. We more particularly focused on penalising the trace of the approximate operator  $T_\nu$ ; however, other type of penalisations may be considered and the study of “optimal” penalisation strategies may deserve further investigations: relations with the spectrum of the initial operator  $T_\mu$ , impact on the sparsity-accuracy trade-off of the induced approximation, etc. The QPs related to the regularised or constrained squared-kernel-discrepancy minimisations problem can also be related to one-class SVMs with respect to the squared-kernel and the initial measure  $\mu$ ; a deeper analysis of this connection may potentially lead to interesting developments.

We investigated two numerical strategies for the resolutions of the regularised and constrained problems. The regularisation-path approach can be used to explore the range of very sparse solutions, with the interest of leading to exact solutions (up to precision errors); the vertex-exchange-based strategy permits the exploration of a wider range of solutions and offer a numerically efficient approach to build approximate solutions. We also proposed two greedy heuristics to enforce the sparsity of a measure obtained from squared-kernel-discrepancy minimisation; the strong and weak-pairwise merging algorithms appear as relatively efficient and very useful in practice. The study of efficient numerical strategies for solving problems related to squared-kernel-discrepancy minimisations is however still obviously largely open and is of great importance for applications.

As illustrated in Section 9, the proposed approach can be used to tackle relatively large-scale

problems while leading to accurate sparse approximations of the main eigenpairs of the considered operator  $T_\mu$  in a reasonable amount of time (and without necessarily resorting to very powerful computing hardwares). In practice, the most computationally demanding steps are the computation of the distortion term  $\mathbf{S}\boldsymbol{\omega}$  and of the Upsilon-test; the computation of  $\mathbf{S}\boldsymbol{\omega}$  is mandatory and can not be avoided, while the Upsilon-test is facultative, although very useful (notice that performing the Upsilon-test for a large number of approximate eigendirections is significantly more demanding than the computation of the distortion term). When dealing with sparse measures, these operations are indeed the only ones with computational complexity scaling as  $\mathcal{O}(N^2)$ ; however, as already mentioned, they can be easily parallelised since they consist in simple kernelised matrix-vector products.

In this work, we have mainly focused on the quadrature-sparsification problem (i.e., the initial measure  $\mu$  is discrete, and the support of  $\nu$  is included in the support of  $\mu$ ); however, as shown in Section 2, the notion of squared-kernel discrepancy appears as a natural tool in the general framework of Nyström approximation of integral operators defined from symmetric positive-semidefinite kernels. From a theoretical point of view, the study of “squared-kernel-discrepancy-optimal” quadrature for Nyström approximation of operators related to general measures  $\mu$  may lead to interesting developments; for instance, the characterisation of discreteness-inducing penalisations (i.e., penalisations leading to discrete optimal measures  $\nu^*$  with finite support) appears of great interest. More generally, the study of kernel-discrepancy-based strategies (not necessarily related to squared-kernels, see for instance [3]) for the optimal sparsification of other types of kernel-based models (like for instance SVM-classification models) may potentially be a fruitful research direction.

**Acknowledgements.** The research leading to these results has received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007-2013) / ERC AdG A-DATADRIVE-B (290923). This paper reflects only the authors’ views, the Union is not liable for any use that may be made of the contained information – Research Council KUL: GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC), BIL12/11T, PhD/Postdoc grants – Flemish Government: FWO: PhD/Postdoc grants, projects: G.0377.12 (Structured systems), G.088114N (Tensor based data similarity); IWT: PhD/Postdoc grants, projects: SBO POM (100031); iMinds Medical Information Technologies SBO 2014 – Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimisation, 2012-2017).

### Appendix A. Proofs.

*Proof of Lemma 2.1.* Consider an o.n.b.  $\{h_k | k \in \mathbb{I}\}$  of  $\mathcal{H}$ . By definition of  $T_\mu$  and  $T_\nu$ , for all  $k \in \mathbb{I}$ , we have

$$\begin{aligned} (T_\mu[h_k] | T_\nu[h_k])_{\mathcal{H}} &= (h_k | T_\nu[h_k])_{L^2(\mu)} = (T_\mu[h_k] | h_k)_{L^2(\nu)} \\ &= \int_{\mathcal{X} \times \mathcal{X}} \mathbf{K}(x, t) h_k(x) h_k(t) d\mu(x) d\nu(t). \end{aligned}$$

For  $x$  and  $t \in \mathcal{X}$ , we have  $K(x, t) = \sum_{k \in \mathbb{I}} h_k(x) h_k(t)$ , and we thus obtain

$$(T_\mu | T_\nu)_{\text{HS}, \mathcal{H}} = \int_{\mathcal{X} \times \mathcal{X}} (K(x, t))^2 d\mu(x) d\nu(t) = \|K\|_{L^2(\mu \otimes \nu)}^2.$$

From Section 2.1, we also have  $\|K\|_{L^2(\mu \otimes \mu)}^2 = \|K_\mu\|_{L^2(\mu \otimes \mu)}^2 = \sum_{k \in \mathbb{I}_\mu^+} \lambda_k^2$ .  $\square$

*Proof of Theorem 2.1.* The proof directly follows from the properties discussed in Sections 2.1 and 2.3. In particular, (2.1) is obtained by considering the o.n.b.  $\{\sqrt{\lambda_k} \varphi_k | k \in \mathbb{I}_\mu^+\}$  of  $\mathcal{H}_\mu$  defined by  $T_\mu$  while remarking that  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$  implies  $T_\mu[h] = T_\nu[h] = 0$  for all  $h \in \mathcal{H}_{0, \mu}$ . By expanding  $T_\nu[\varphi_k]$  thanks to  $P_{\mathcal{H}_\nu}[\varphi_k] = \sum_{l \in \mathbb{I}_\nu^+} \vartheta_l \psi_l(\psi_l | \varphi_k)_{\mathcal{H}}$  (orthogonal projection, in  $\mathcal{H}$ , of  $\varphi_k$  onto  $\mathcal{H}_\nu$ ), we obtain (2.2). Finally, the last inequality directly follows from  $\|h\|_{L^2(\mu)}^2 \leq \tau_\mu \|h\|_{\mathcal{H}}^2$  for all  $h \in \mathcal{H}$ .  $\square$

*Proof of Lemma 3.1.* If  $\boldsymbol{\omega} = \mathbf{1}$  and the components of  $\mathbf{v}$  are all 0 or 1, then the components of  $\boldsymbol{\omega} - \mathbf{v}$  are also all 0 or 1. Introducing the index set  $I = \{k | v_k = 0\} = \{k | \omega_k - v_k = 1\}$ , we obtain  $(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) = \sum_{(i, j) \in I \times I} K^2(x_i, x_j)$ , which is equal to  $\|\mathbf{K} - \mathbf{V}\mathbf{K}\mathbf{V}\|_F^2$ .  $\square$

*Proof of Lemma 3.2.* Assertion (i) follows directly from  $D_{K^2}(\mu, \mu) = 0$  (since  $D_{K^2}(\mu, \nu) \geq 0$ ). From the first order optimality condition, for  $\alpha \geq 0$ , a feasible  $\mathbf{v}_\alpha^*$  is solution to (3.1) if and only if, for any feasible  $\mathbf{v}$ , we have  $(\mathbf{v} - \mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) \geq 0$ . Then, considering  $\mathbf{v}_\alpha^* = 0$  leads to assertion (ii), and since  $\boldsymbol{\omega}$  is by assumption feasible for (3.1), assertion (iii) is obtained by taking  $\mathbf{v} = \boldsymbol{\omega}$ . To prove (iv), we first remark that for all  $\theta \in [0, 1]$ , we have

$$D_\alpha(\theta \mathbf{v}_\alpha^* + (1 - \theta) \tilde{\mathbf{v}}_\alpha^*) = D_\alpha(\tilde{\mathbf{v}}_\alpha^*) + \theta (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \nabla D_\alpha(\tilde{\mathbf{v}}_\alpha^*) + \theta^2 \frac{1}{2} (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \mathbf{S} (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*).$$

Since  $D_\alpha(\tilde{\mathbf{v}}_\alpha^*) = D_\alpha(\theta \mathbf{v}_\alpha^* + (1 - \theta) \tilde{\mathbf{v}}_\alpha^*)$ , we necessarily have  $(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \mathbf{S} (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*) = 0$ . From the Cauchy-Schwarz inequality for the bilinear form related to  $\mathbf{S}$ , we therefore have  $(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \mathbf{S} (\tilde{\mathbf{v}}_\alpha^* - \boldsymbol{\omega}) = 0$ , and this completes the proof since we also have  $(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \nabla D_\alpha(\tilde{\mathbf{v}}_\alpha^*) = 0$ . Assertion (v) follows from the expansion, for all  $\mathbf{v} \geq 0$ ,

$$D_\alpha(\mathbf{v}) = D_\alpha(\mathbf{v} - v_k \mathbf{e}_k) + [\boldsymbol{\alpha} \mathbf{d} - \mathbf{S} \boldsymbol{\omega}]_k v_k + \sum_{i \neq k} v_i v_k \mathbf{S}_{i,k} + \frac{1}{2} v_k^2 \mathbf{S}_{k,k},$$

where  $\mathbf{e}_k$  stands for the  $k$ -th element of the canonical basis of  $\mathbb{R}^N$  (so that  $[\mathbf{v} - v_k \mathbf{e}_k]_k = 0$ ). Since all the entries of  $\mathbf{S}$  are non-negative, if  $[\boldsymbol{\alpha} \mathbf{d} - \mathbf{S} \boldsymbol{\omega}]_k \geq 0$  and  $v_k > 0$ , since  $\mathbf{S}_{k,k} > 0$ , we therefore have  $D(\mathbf{v} - v_k \mathbf{e}_k) < D_\alpha(\mathbf{v})$ . To obtain assertion (vi), consider  $\alpha_1 < \alpha_2$  and let  $\mathbf{v}_{\alpha_1}^*$  and  $\mathbf{v}_{\alpha_2}^*$  be solutions to problem (3.1) with  $\alpha = \alpha_1$  and  $\alpha = \alpha_2$  respectively. We have  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \leq \alpha_1 \mathbf{d}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*)$  and  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \geq \alpha_2 \mathbf{d}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*)$ , so that, necessarily,  $\mathbf{d}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*) \leq 0$ , and therefore  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \leq 0$ .  $\square$

*Proof of Lemma 3.3.* If  $\mathbf{v}_\alpha^*$  is a solution to problem (3.1) with  $\alpha \geq 0$ , then by definition,  $\mathbf{v}_\alpha^*$  in particular minimises  $D(\cdot)$  over the set of all  $\mathbf{v} \geq 0$  such that  $\mathbf{d}^T \mathbf{v} = \mathbf{d}^T \mathbf{v}_\alpha^*$ , and  $\mathbf{v}_\alpha^*$  is therefore a solution to problem (3.2) with  $\varkappa = \mathbf{d}^T \mathbf{v}_\alpha^*$ .

The condition  $\varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$  follows directly from Lemma 3.2-(iii): a solution  $\mathbf{v}_\alpha^*$  to problem (3.1) indeed necessarily satisfies  $\mathbf{d}^T \mathbf{v}_\alpha^* \leq \mathbf{d}^T \boldsymbol{\omega}$ . For  $\varkappa = 0$ , we have  $\mathbf{v}_\alpha^* = 0$  and the result follows from Lemma 3.2-(ii). For  $0 < \varkappa \leq \mathbf{d}^T \boldsymbol{\omega}$ , in order to be a solution to problem (3.1), a solution  $\mathbf{v}_\alpha^*$  to problem (3.2) must satisfy (first order optimality condition),  $(\mathbf{v} - \mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) \geq 0$  for all  $\mathbf{v} \geq 0$ . In particular, for all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ , the constraint  $v_k \geq 0$  is not active; therefore, for all  $\theta \geq 0$ , the parameter  $\alpha$  must satisfy  $(\theta \mathbf{e}_k - \mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) = 0$  where  $\mathbf{e}_k$  stands for the  $k$ -th element of the canonical basis of  $\mathbb{R}^N$  (so that  $[\mathbf{v}_\alpha^*]_k = \mathbf{e}_k^T \mathbf{v}_\alpha^*$ ). Considering  $\theta = 0$  directly leads to the expected result (and one can also easily verify that the obtained value of  $\alpha$  does not depend on  $\theta \geq 0$ ). The last assertion follows directly from Lemma 3.2-(vi) and the correspondence between solutions to problems (3.1) and (3.2).  $\square$

*Proof of Lemma 4.1.* The proof follows directly from the expansion of the involved squared norms. We indeed have  $\|T_\mu[\hat{\varphi}_l]\|_{L^2(\mu)}^2 = \hat{\lambda}_l^2$ ,  $\|\hat{\lambda}_l \hat{\varphi}_l\|_{L^2(\mu)}^2 = \hat{\lambda}_l^2$ , and  $(\hat{\lambda}_l \hat{\varphi}_l | T_\mu[\hat{\varphi}_l])_{L^2(\mu)} = \hat{\lambda}_l^2 \Upsilon_l$ ; and  $\|T_\mu[\hat{\varphi}_l]\|_{\mathcal{H}}^2 = \hat{\lambda}_l \Upsilon_l$ ,  $\|\hat{\lambda}_l \hat{\varphi}_l\|_{\mathcal{H}}^2 = \hat{\lambda}_l \|\hat{\lambda}_l^{1/2} \hat{\varphi}_l\|_{\mathcal{H}}^2$ , and  $(\hat{\lambda}_l \hat{\varphi}_l | T_\mu[\hat{\varphi}_l])_{\mathcal{H}} = \hat{\lambda}_l$ .  $\square$

*Proof of Lemma 5.1.* Define the closed linear subspace  $\mathcal{F}_S = \text{span}\{K_{x_k}^2\}_{k=1}^N$  of  $\mathcal{F}$  and let  $\mathcal{F}_0 = \mathcal{F}_S^\perp$  be its orthogonal; by definition,  $g_\omega \in \mathcal{F}_S$ . For any  $g_S \in \mathcal{F}_S$  and  $g_0 \in \mathcal{F}_0$ , we have

$$\frac{1}{2} \|g_S\|_{\mathcal{F}}^2 + (g_S | g_\omega)_{\mathcal{F}} \leq \frac{1}{2} \|g_S + g_0\|_{\mathcal{F}}^2 + (g_S + g_0 | g_\omega)_{\mathcal{F}} = \frac{1}{2} \|g_S\|_{\mathcal{F}}^2 + (g_S | g_\omega)_{\mathcal{F}} + \frac{1}{2} \|g_0\|_{\mathcal{F}}^2.$$

In addition, for any  $k \in \{1, \dots, N\}$ , we have  $g_0(x_k) = 0$ , so that necessarily  $g_\alpha^* \in \mathcal{F}_S$  (representer theorem), i.e, there exists  $\boldsymbol{\beta}^* = (\beta_1^*, \dots, \beta_N^*)^T \in \mathbb{R}^N$  such that  $g_\alpha^* = \sum_{k=1}^N \beta_k^* K_{x_k}^2$ . Restricting problem (5.1) to  $\mathcal{F}_S$  then yields, for  $\boldsymbol{\beta} \in \mathbb{R}^N$ ,

$$\text{minimise}_{\boldsymbol{\beta}} \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} \text{ subject to } \mathbf{S} \boldsymbol{\beta} \geq -\boldsymbol{\alpha} \mathbf{d}. \quad (\text{A.1})$$

We then introduce the Lagrangian function, for  $\mathbf{v} \in \mathbb{R}^N$  with  $\mathbf{v} \geq 0$  (dual feasibility conditions),

$$\mathcal{L}(\boldsymbol{\beta}, \mathbf{v}) = \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} - \mathbf{v}^T [\mathbf{S} \boldsymbol{\beta} + \boldsymbol{\alpha} \mathbf{d}].$$

The primal optimality conditions give  $\mathbf{S}\boldsymbol{\beta} = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega})$ , leading to the Lagrange dual problem (3.1) (written as a minimisation problem). If  $\mathbf{v}_\alpha^*$  is a solution to (3.2), then a solution  $\boldsymbol{\beta}^*$  to (A.1) needs to satisfy  $\mathbf{S}\boldsymbol{\beta}^* = \mathbf{S}(\mathbf{v}_\alpha^* - \boldsymbol{\omega})$ , so that we can in particular consider  $\boldsymbol{\beta}^* = \mathbf{v}_\alpha^* - \boldsymbol{\omega}$ . Notice that when  $\mathbf{S}$  is non-invertible, other choices for  $\boldsymbol{\beta}^*$  exist since for any  $\boldsymbol{\varepsilon} \in \mathbb{R}^N$  such that  $\mathbf{S}\boldsymbol{\varepsilon} = 0$ , we have  $\mathbf{S}(\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}) = \mathbf{S}\boldsymbol{\beta}^*$ ; but the obtained solution  $g_\alpha^* \in \mathcal{F}_S$  does not depend on such a  $\boldsymbol{\varepsilon}$  (see also Remark 3.2). The equality  $g_\alpha^*(x_k) = -\alpha d_k$  for all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\alpha^*]_k > 0$  is consequence of the complementary slackness condition  $(\mathbf{v}_\alpha^*)^T [\mathbf{S}(\mathbf{v}_\alpha^* - \boldsymbol{\omega}) + \alpha \mathbf{d}] = 0$ .  $\square$

*Proof of Lemma 5.2.* We follow the same reasoning than in the proof of Lemma 5.1. By restricting problem (5.3) to  $\mathcal{F}_S$ , we obtain, for  $\boldsymbol{\beta} \in \mathbb{R}^N$ ,

$$\underset{\boldsymbol{\beta}, \gamma}{\text{minimise}} \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} - \gamma \text{ subject to } \mathbf{S} \boldsymbol{\beta} \geq \gamma \mathbf{d} / \boldsymbol{\varkappa}. \quad (\text{A.2})$$

The underlying Lagrangian function is then given by, for  $\mathbf{v} \in \mathbb{R}^N$  with  $\mathbf{v} \geq 0$  (dual feasibility conditions),

$$\mathcal{L}(\boldsymbol{\beta}, \gamma, \mathbf{v}) = \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} - \gamma - \mathbf{v}^T [\mathbf{S} \boldsymbol{\beta} - \gamma \mathbf{d} / \boldsymbol{\varkappa}].$$

The primal optimality conditions give  $\mathbf{S} \boldsymbol{\beta} = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega})$  and  $\mathbf{d}^T \mathbf{v} = \boldsymbol{\varkappa}$ , leading to the Lagrange dual problem (3.2). If  $\mathbf{v}_x^*$  is a solution to (3.2), then a solution  $\boldsymbol{\beta}^*$  to (A.2) needs to satisfy  $\mathbf{S} \boldsymbol{\beta}^* = \mathbf{S}(\mathbf{v}_x^* - \boldsymbol{\omega})$ , so that we can in particular consider  $\boldsymbol{\beta}^* = \mathbf{v}_x^* - \boldsymbol{\omega}$ . The expression of  $\gamma_x^*$  follows from the complementary slackness condition  $(\mathbf{v}_x^*)^T [\mathbf{S}(\mathbf{v}_x^* - \boldsymbol{\omega}) - \gamma_x^* \mathbf{d} / \boldsymbol{\varkappa}] = 0$ . The equality  $g_x^*(x_k) = \gamma_x^* d_k / \boldsymbol{\varkappa}$  for all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_x^*]_k > 0$  is also consequence of the complementary slackness condition.  $\square$

*Proof of Lemma 6.1.* Let  $\mathbf{v}_\alpha = \theta \mathbf{v}_{\alpha_1}^* + (1 - \theta) \mathbf{v}_{\alpha_2}^*$  and define  $J = J_{\alpha_1} = J_{\alpha_2}$ ; we have

$$\mathbf{S}_{J,J} [\mathbf{v}_\alpha]_J = \mathbf{S}_{J,J} [\theta \mathbf{v}_{\alpha_1}^* + (1 - \theta) \mathbf{v}_{\alpha_2}^*]_J = [\mathbf{S} \boldsymbol{\omega}]_J - \alpha \mathbf{d}_J,$$

so that  $[\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \alpha \mathbf{d}]_J = 0$ , and in the same way,

$$[\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \alpha \mathbf{d}]_{J^c} = \theta [\mathbf{S}(\mathbf{v}_{\alpha_1}^* - \boldsymbol{\omega}) + \alpha_1 \mathbf{d}]_{J^c} + (1 - \theta) [\mathbf{S}(\mathbf{v}_{\alpha_2}^* - \boldsymbol{\omega}) + \alpha_2 \mathbf{d}]_{J^c} > 0.$$

By construction,  $\mathbf{v}_\alpha \geq 0$  and in addition, if  $k$  is such that  $[\mathbf{v}_\alpha]_k > 0$ , then  $k \in J$  (since these conditions are verified by both  $\mathbf{v}_{\alpha_1}^*$  and  $\mathbf{v}_{\alpha_2}^*$ ). We therefore have  $\mathbf{v}_\alpha^T (\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \alpha \mathbf{d}) = 0$ , so that for all  $\mathbf{v} \geq 0$ , the optimality condition  $(\mathbf{v} - \mathbf{v}_\alpha)^T \nabla D_\alpha(\mathbf{v}_\alpha) \geq 0$  holds, i.e.,  $\mathbf{v}_\alpha$  is a solution to (3.1), and  $J_\alpha = J$ .  $\square$

## REFERENCES

- [1] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science, 2011.
- [2] Christos Boutsidis, Petros Drineas, and Malik Magdon-Ismael. Near-optimal column-based matrix reconstruction. *SIAM Journal on Computing*, 43(2):687–717, 2014.
- [3] Steve B. Damelin. A walk through energy, discrepancy, numerical integration and group invariant measures on measurable subsets of Euclidean space. *Numerical Algorithms*, 48(1-3):213–235, 2008.
- [4] Petros Drineas and Michael W. Mahoney. On the Nyström method for approximating a Gram matrix for improved kernel-based learning. *The Journal of Machine Learning Research*, 6:2153–2175, 2005.
- [5] Bertrand Gauthier and Luc Pronzato. Spectral approximation of the IMSE criterion for optimal designs in kernel-based interpolation models. *SIAM/ASA Journal on Uncertainty Quantification*, 2:805–825, 2014.
- [6] Alex Gittens and Michael W. Mahoney. Revisiting the Nyström method for improved large-scale machine learning. *Journal of Machine Learning Research*, 17:1–65, 2016.
- [7] Jacques Guélat and Patrice Marcotte. Some comments on Wolfe’s “away step”. *Mathematical Programming*, 35(1):110–119, 1986.
- [8] Wolfgang Hackbusch. *Integral equations: theory and numerical treatment*, volume 120. Birkhäuser, 2012.
- [9] Trevor Hastie, Saharon Rosset, Robert Tibshirani, and Ji Zhu. The entire regularization path for the Support Vector Machine. *Journal of Machine Learning Research*, 5:1391–1415, 2004.
- [10] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical Learning with Sparsity: the Lasso and Generalizations*. CRC Press, 2015.



- [11] Sanjiv Kumar, Mehryar Mohri, and Ameet Talwalkar. Sampling methods for the Nyström method. *Journal of Machine Learning Research*, 13:981–1006, 2012.
- [12] Moshe Lichman. UCI Machine Learning Repository, 2013.
- [13] Michael R. Osborne, Brett Presnell, and Berwin A. Turlach. A new approach to variable selection in least squares problems. *IMA journal of numerical analysis*, 20(3):389–403, 2000.
- [14] Mee Young Park and Trevor Hastie. L1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(4):659–677, 2007.
- [15] Luc Pronzato and Andrej Pázman. *Design of Experiments in Nonlinear Models*. Springer, 2013.
- [16] Bernhard Schölkopf, John C. Platt, John Shawe-Taylor, Alexander J. Smola, and Robert C. Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, 2001.
- [17] Laurent Schwartz. *Analyse Hilbertienne*. Hermann, 1978.
- [18] Steve Smale and Ding-Xuan Zhou. Geometry on probability spaces. *Constructive Approximation*, 30(3):311–323, 2009.
- [19] Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer, 2008.
- [20] Shusen Wang and Zhihua Zhang. Improving CUR matrix decomposition and the Nyström approximation via adaptive sampling. *The Journal of Machine Learning Research*, 14:2729–2769, 2013.