



**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. 2017. hal-01416786v2

**HAL Id: hal-01416786**

**<https://hal.science/hal-01416786v2>**

Preprint submitted on 30 Mar 2017 (v2), 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. We more specifically focus on computing sparse quadratures with support included in fixed finite sets of points (quadrature-sparsification), this framework being of particular interest since it encompasses the Nyström approximation of kernel-matrices. For a given kernel, the accuracy of a quadrature approximation is assessed through the squared Hilbert-Schmidt norm (for operators on the underlying reproducing kernel Hilbert space) of the difference between the integral operators related to the initial and approximate measures; by analogy with the notion of kernel discrepancy, we refer to the underlying criterion as the squared-kernel discrepancy between two measures. Sparsity of the approximate quadrature is achieved through the introduction of an  $\ell^1$ -type penalisation, and the computation of a penalised squared-kernel-discrepancy-optimal approximation thus consists in 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. The quadratic programs related to the regularised and constrained problems can in particular be interpreted as the Lagrange dual formulations of distorted one-class support-vector machines related to the squared kernel and the initial 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 sparse representations 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.** Computing 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. Such a situation is for instance encountered each time a pointwise quadrature is used to approximate an integral operator, or more directly when one aims at computing the eigendecomposition of a kernel-matrix (or of any Gram matrix). In the non-sparse case and for the direct approach, the amount of computations required to perform the diagonalisation of such a matrix scales as  $\mathcal{O}(N^3)$ , and thus becomes quickly intractable when  $N$  is large (not to mention issues related to the storage of large-scale matrices).

In the framework of matrix approximation, 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  $N$ ; by analogy with the integral operators framework (see, e.g., [9]), this approach is usually referred to as the Nyström method. Naturally, the choice of the considered submatrix has a strong impact on the quality of the induced approximation, raising questions relative to the selection of an appropriate submatrix, see for instance [4, 21, 12, 23, 2, 7] and references therein. Translated in the integral-operator framework, this operation consists in defining a sparse approximation of an initial discrete measure, the support of the approximate measure being included in the support of the initial one. More generally, we refer to the problem consisting in computing a sparse quadrature while enforcing the support of this quadrature to be included in a fixed finite set of points as *quadrature-sparsification*.

Following for instance [21], for a given kernel  $K(\cdot, \cdot)$  and an initial measure  $\mu$  (defining an integral operator  $T_\mu$ , see Section 2 for a detailed discussion), we assess the accuracy of an approximate measure  $\nu$  (defining an operator  $T_\nu$ ) by considering the squared Hilbert-Schmidt norm of the difference between the integral operators  $T_\mu$  and  $T_\nu$  viewed as operators on the underlying reproducing kernel Hilbert space (RKHS, see for instance [1]). By analogy with the notion of kernel discrepancy, see [3] and Appendix A, we refer to the underlying criterion as the *squared-kernel discrepancy* between the

---

\*gauthierb@cardiff.ac.uk (corresponding author)

‡johan.suykens@esat.kuleuven.be

§Cardiff University, School of Mathematics

Senghennydd Road, Cardiff, CF24 4AG, United Kingdom

†KU Leuven, ESAT-STADIUS Centre for Dynamical Systems, Signal Processing and Data Analytics  
Kasteelpark Arenberg 10, B-3001 Leuven, Belgium

measures  $\mu$  and  $\nu$ , denoted by  $D_{K^2}(\mu, \nu)$ .

For a given  $n$ , the search of a measure  $\nu_n^*$  supported by  $n$  points and such that  $D_{K^2}(\mu, \nu_n^*)$  is minimal, is in general a difficult non-convex minimisation problem. As an alternative, we assume that the support of the quadrature  $\nu$  is included in a fixed finite set of  $N$  points (quadrature-sparsification; this framework in particular encompasses the Nyström approximation of kernel-matrices). The function  $\nu \mapsto D_{K^2}(\mu, \nu)$  then consists in a convex quadratic function on  $\mathbb{R}^N$ , and sparsity of the approximate quadrature can be induced by the introduction of an  $\ell^1$ -type penalisation of the squared-kernel discrepancy (the sample size  $n$  is thus not fixed a priori). The penalisation can be introduced under the form of a regularisation term or of a constraint (both formulations being equivalent), and the computation of a penalised squared-kernel-discrepancy-optimal quadrature with support included in  $S$  then consists in a convex quadratic minimisation problem (see Section 3). The quadratic programs (QPs) related to the regularised and constrained squared-kernel-discrepancy minimisation problems can in particular be interpreted as the Lagrange dual formulations of distorted one-class support-vector machines related to the squared kernel and the initial measure, as discussed in Section 4.

We present a careful analysis of the approach and present two different numerical strategies for solving penalised squared-kernel-discrepancy minimisation problems. The first strategy, described in Section 5, is a sequential direct method based on the notion of regularisation path (see for instance [16]); this method appears as an efficient way to explore the range of very sparse solutions to the regularised or constrained problems. The second strategy, see Section 6, applies more specifically to the constrained formulation and consists in using a vertex-exchange algorithm (see, e.g., [17, Chap. 9]) to build an approximate solution to the underlying QP. Importantly, the proposed strategies do not require the storage of any large objects, so that large-scale problems can be considered (more particularly since we are interested in the range of sparse solutions). In addition, to enhance the sparsity of a given discrete measure  $\nu$  while trying to keep the squared-kernel discrepancy as low as possible, we describe two greedy heuristics based on iterative pairwise-component merging (Section 6.2).

We also further investigate the problem related to the computation of an approximate eigendecomposition of the initial operator  $T_\mu$  from the eigendecomposition of an approximate operator  $T_\nu$ . Based on geometrical considerations, we introduce four different approximations (referred to as *approximate eigenvalues*) of the eigenvalue related to a given approximate eigendirection of  $T_\mu$  (i.e., the eigendirections of the approximate operator  $T_\nu$ ), and the accuracy of a given approximate eigendirection can then in particular be assessed through the comparison of these different approximate eigenvalues, as detailed in Section 7. A discussion relative to the choice of relevant penalisation weights (since a weighted- $\ell^1$ -type penalisation is considered) is proposed Section 8. Finally, Sections 9 (two-dimension toy example) and 10 (large-scale problems) are devoted to numerical experiments, and Section 11 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 C.

**2. Notations, recalls and theoretical motivations.** This section is devoted to recalls related to integral operators defined from a symmetric positive-semidefinite kernel. We consider a general space  $\mathcal{X}$  and a symmetric and positive-semidefinite kernel  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ ; we denote by  $\mathcal{H}$  the underlying reproducing RKHS of real-valued functions on  $\mathcal{X}$  (see for instance [1]). We assume that  $\mathcal{H}$  is a separable Hilbert space.

**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 [22, 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

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

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

of  $K(\cdot, \cdot)$  and the Cauchy-Schwarz inequality for the inner product of  $\mathcal{H}$ )

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

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 can thus define the symmetric and positive-semidefinite integral operator  $T_\mu$  on  $L^2(\mu)$ , given by, for  $f \in L^2(\mu)$  and  $x \in \mathcal{X}$ ,

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

We have  $T_\mu[f] \in \mathcal{H}$  for all  $f \in L^2(\mu)$  and, for  $h \in \mathcal{H}$ ,  $(h|T_\mu[f])_{\mathcal{H}} = (h|f)_{L^2(\mu)}$ , see for instance [5, 6] for more details.

We introduce the closed linear subspaces  $\mathcal{H}_{0\mu} = \{h \in \mathcal{H} | \|h\|_{L^2(\mu)} = 0\}$  and  $\mathcal{H}_\mu = \mathcal{H}_{0\mu}^{\perp\mathcal{H}}$  (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  $\{\lambda_k | k \in \mathbb{N}^+\}$  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  $\{\sqrt{\lambda_k} \varphi_k | k \in \mathbb{N}^+\}$  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  $K_{0\mu}(\cdot, \cdot) = K(\cdot, \cdot) - K_\mu(\cdot, \cdot)$  is the reproducing kernel of  $\mathcal{H}_{0\mu}$ ; We also recall that  $\tau_\mu = \sum_{k \in \mathbb{N}^+} \lambda_k$  is the trace of the integral operator  $T_\mu$  on  $L^2(\mu)$ .

**2.2. Discrete measures and kernel-matrices .** Let  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  be a discrete measure supported by  $\mathcal{S} = \{x_k\}_{k=1}^N$ , with  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_N)^T \in \mathbb{R}^N$ ,  $\omega_k > 0$  (in what follows, we 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.,  $\mathbf{K} \mathbf{W} = \mathbf{P} \boldsymbol{\Lambda} \mathbf{P}^{-1}$  with  $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_N)$  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  $\mathbf{K} = \mathbf{P} \boldsymbol{\Lambda} \mathbf{P}^T$ , 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; however, 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  $\mathbf{W}^{1/2} \mathbf{v}_k$ . Numerically, 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  $\mathcal{H}$  (see, e.g., [21]); 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.  $\{h_j | j \in \mathbb{I}\}$  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

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

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., [20]. The underlying Hilbert-Schmidt norm (for operators on  $\mathcal{H}$ ) is given by

$$\|T_\mu\|_{\text{HS}(\mathcal{H})}^2 = (T_\mu | T_\mu)_{\text{HS}(\mathcal{H})} = \sum_{j \in \mathbb{I}} \|T_\mu[h_j]\|_{\mathcal{H}}^2.$$

**Definition 2.1.** *The squared-kernel discrepancy between  $\mu$  and  $\nu \in \mathcal{T}(K)$  is defined as*

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

**Lemma 2.1.** *For  $\mu$  and  $\nu \in \mathcal{T}(K)$ , we have  $(T_\mu | T_\nu)_{\text{HS}(\mathcal{H})} = \|K\|_{L^2(\mu \otimes \nu)}^2$ , so that*

$$D_{K^2}(\mu, \nu) = \|K\|_{L^2(\mu \otimes \mu)}^2 + \|K\|_{L^2(\nu \otimes \nu)}^2 - 2\|K\|_{L^2(\mu \otimes \nu)}^2,$$

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

In particular, note that  $\|K\|_{L^2(\mu \otimes \nu)}^2 \leq \tau_\mu \tau_\nu$ , and that  $\|T_\mu\|_{\text{HS}(\mathcal{H})}^2 = \sum_{k \in \mathbb{I}_\mu^+} \lambda_k^2$ , where  $\{\lambda_k | k \in \mathbb{I}_\mu^+\}$  is the set of all strictly positive eigenvalues of  $T_\mu$ . By definition, we always have  $D_{K^2}(\mu, \nu) \geq 0$ , and  $D_{K^2}(\mu, \mu) = 0$ . We can 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  $\|K\|_{L^2(\mu \otimes \nu)}^2 = 0$ .

The terminology ‘‘squared-kernel discrepancy’’ is motivated by the analogy with the notion of ‘‘kernel discrepancy’’ discussed for instance in [3] (see Appendix A). Interestingly, the kernel discrepancy is related to approximate integration of functions in the RKHS  $\mathcal{H}$ , while the squared-kernel discrepancy is related to the approximation of integral operators defined from the reproducing kernel  $K(\cdot, \cdot)$  of  $\mathcal{H}$ . The squared-kernel discrepancy is also related to approximate integration of functions in the RKHS associated with the *squared kernel*; in what follows we denote by  $\mathcal{G}$  the RKHS associated with  $K^2(\cdot, \cdot) = (K(\cdot, \cdot))^2$ . Following Appendix A, for all  $\mu \in \mathcal{T}(K)$  (notice that  $\mathcal{T}(K) = \mathcal{T}(K^2)$ ), we can introduce  $g_\mu(x) = \int_{\mathcal{X}} K^2(x, t) d\mu(t)$ , with  $x \in \mathcal{X}$ . We have  $g_\mu \in \mathcal{G}$  and finally

$$D_{K^2}(\mu, \nu) = \|g_\mu - g_\nu\|_{\mathcal{G}}^2. \quad (2.1)$$

The following Lemma 2.2 further highlights the connection between the squared-kernel discrepancy and the error induced by the approximation of  $T_\mu$  (the initial integral operator) by  $T_\nu$  (the approximate operator).

**Lemma 2.2.** *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  $T_\mu$ , we have*

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

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

In view of (2.2), for a fixed measure  $\mu$ , by minimising the function  $\nu \mapsto D_{K^2}(\mu, \nu)$  under the condition  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$ , we minimise, 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 rule of penalisation weights). When  $D_{K^2}(\mu, \nu)$  is small, we can therefore expect the main eigendirections of  $T_\nu$  to be accurate approximations of the main eigendirections of  $T_\mu$  (i.e., the ones related to the largest eigenvalues).

Since  $D_{K^2}(\mu, \mu) = 0$  (i.e., ‘‘the best approximation of  $T_\mu$  is  $T_\mu$  itself’’), the unconstrained minimisation of  $\nu \mapsto D_{K^2}(\mu, \nu)$  on  $\mathcal{T}(K)$  is of no interest. In the framework of sparse pointwise

quadrature approximation, we aim at obtaining a discrete measure  $\nu$  supported by a relatively small number of points (in order to be able to compute the eigendecomposition of  $T_\nu$ ) and related to an as low as possible value of  $D_{K^2}(\mu, \nu)$ . However, for a given  $n \in \mathbb{N}^*$ , the search of an optimal discrete measure  $\nu_n^*$  such that  $D_{K^2}(\mu, \nu_n^*)$  is minimal among all measures  $\nu_n$  supported by  $n$  points is in general a difficult (i.e, usually non-convex) optimisation problem on  $(\mathcal{X} \times \mathbb{R}_+)^n$ . To avoid this difficulty, in the following Section 3, we restrict the problem to measures  $\nu$  with support included in a fixed finite set of points  $S = \{x_k\}_{k=1}^N$  (with, in practice,  $N$  large), and instead of fixing a priori the number  $n$  of support points, we induce sparsity through the introduction of an  $\ell^1$ -type penalisation. We refer to the underlying problem as *quadrature-sparsification*.

**3. Optimal quadrature-sparsification as quadratic programming.** We consider a fixed set of points  $S$ , and we restrict our study to discrete measures  $\nu$  with support included in  $S$ . We more particularly focus on the case where  $S$  is the support of an initial discrete measure  $\mu$ .

**3.1. Preliminary discussion.** Let  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$ , with  $\mathbf{v} = (v_1, \dots, v_N)^T \geq 0$  (that is  $v_k \geq 0$  for all  $k$ ), and consider a general fixed measure  $\mu \in \mathcal{T}(K)$ ; we have

$$\|K\|_{L^2(\nu \otimes \nu)}^2 = \mathbf{v}^T \mathbf{S} \mathbf{v} \text{ and } \|K\|_{L^2(\mu \otimes \nu)}^2 = \mathbf{g}_\mu^T \mathbf{v},$$

where  $\mathbf{S}$  is the kernel-matrix defined by the squared kernel  $K^2(\cdot, \cdot)$  and the set of points  $S$ , i.e., with  $i, j$  entry  $S_{i,j} = K^2(x_i, x_j) \geq 0$  ( $\mathbf{S}$  is therefore a non-negative, positive-semidefinite, symmetric matrix), and where  $\mathbf{g}_\mu = (g_\mu(x_1), \dots, g_\mu(x_N))^T \in \mathbb{R}^N$ , with  $g_\mu(x_k) = \int_{\mathcal{X}} K^2(x_k, t) d\mu(t)$  (so that  $\mathbf{g}_\mu \geq 0$ ). For such a discrete measure  $\nu$ , we thus obtain

$$D_{K^2}(\mu, \nu) = \|K\|_{L^2(\mu \otimes \mu)}^2 + \mathbf{v}^T \mathbf{S} \mathbf{v} - 2\mathbf{g}_\mu^T \mathbf{v}, \quad (3.1)$$

and  $\nu \mapsto D_{K^2}(\mu, \nu)$  can this way be interpreted as a quadratic function of  $\mathbf{v}$  (i.e., the vector of weights characterising  $\nu$ ). By minimising  $\mathbf{v} \mapsto \mathbf{v}^T \mathbf{S} \mathbf{v} - 2\mathbf{g}_\mu^T \mathbf{v}$  under the constraint  $\mathbf{v} \geq 0$ , we can obtain the best approximation of  $\mu$  (in terms of squared-kernel discrepancy) among all discrete measures supported by  $S$  (note that this in practice requires the knowledge of the vector  $\mathbf{g}_\mu \in \mathbb{R}^N$ , which may be problematic for a general measure  $\mu$ ). We however not consider such a minimisation problem; instead, we introduce an  $\ell^1$ -type penalisation of the squared-kernel discrepancy (to obtain approximate measures  $\nu$  supported by a number of points significantly smaller than  $N$ ).

In order to discard problems relative to the computation of the vector  $\mathbf{g}_\mu$  and for the sake of simplicity, we assume that the measure  $\mu$  is also discrete and supported by  $S$ . More precisely, we consider that  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$ , with  $\boldsymbol{\omega} > 0$  (i.e.,  $\omega_k > 0$  for all  $k$ ). We then in particular have  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$  for all  $\nu \geq 0$  (see Lemma 2.2), and  $\mathbf{g}_\mu = \mathbf{S} \boldsymbol{\omega}$ ; in addition, note that this framework is of particular interest since it corresponds to problems related to kernel-matrix low-rank approximation, see Section 2.2, Remark 3.1, and for instance [4].

Thus, for  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  and  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$  be, with  $\boldsymbol{\omega} > 0$  and  $\mathbf{v} \geq 0$ , we have

$$D_{K^2}(\mu, \nu) = (\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S} (\boldsymbol{\omega} - \mathbf{v}). \quad (3.2)$$

For a fixed discrete measure  $\mu$  supported by  $S$  (i.e.,  $\boldsymbol{\omega} > 0$  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.

**Remark 3.1** (relation with the classical Frobenius-norm-based criterion). In the framework of equation (3.2), 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  $\mathbf{W} = \text{diag}(\boldsymbol{\omega}) = \text{Id}_N$  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 interpreted 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 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  $\mathbf{V}^2 = \mathbf{V} = \mathbf{V}^T$ , selecting a principal submatrix of the matrix  $\mathbf{K}$  can be viewed as considering the matrix  $\mathbf{V}^{1/2}\mathbf{KV}^{1/2} = \mathbf{VKV}$ .

In this case, i.e., if  $\boldsymbol{\omega} = \mathbf{1}$  and if the components of  $\mathbf{v}$  are all 0 or 1, we can easily verify that  $(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) = \|\mathbf{K} - \mathbf{VKV}\|_F^2$ , where  $\|\cdot\|_F$  stands for the Frobenius norm.  $\triangleleft$

**3.2. Regularised squared-kernel-discrepancy minimisation.** We consider the framework of equation (3.2) (i.e., the measures  $\mu$  and  $\nu$  are discrete, with support included in  $S$ ). For a given penalisation vector  $\mathbf{f} = (f_1, \dots, f_N)^T \in \mathbb{R}^N$ , with  $\mathbf{f} > 0$  (see Section 8), and  $\alpha \geq 0$ , we introduce the *regularised squared-kernel-discrepancy minimisation problem*

$$\underset{\mathbf{v}}{\text{minimise}} D_\alpha(\mathbf{v}) = \frac{1}{2}(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) + \alpha \mathbf{f}^T \mathbf{v} \text{ subject to } \mathbf{v} \geq 0. \quad (3.3)$$

Notice that  $D_\alpha(\mathbf{v}) = D(\mathbf{v}) + \alpha \mathbf{f}^T \mathbf{v}$ . In particular, when  $\mathbf{S}$  is invertible,  $D_\alpha(\cdot)$  is strongly convex (and, in this case, a solution to (3.3) is thus necessarily unique). We also recall that, for a given  $\alpha$ , the set of solutions to (3.3) is convex. The gradient of  $D_\alpha$  at  $\mathbf{v} \in \mathbb{R}^N$  is given by  $\nabla D_\alpha(\mathbf{v}) = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega}) + \alpha \mathbf{f}$ . The following Lemma 3.1 recalls some simple properties verified by the solutions to problem (3.3).

**Lemma 3.1.** *Denote by  $\mathbf{v}_\alpha^*$  a solution to (3.3) with  $\alpha \geq 0$ , we have:*

- (i) for  $\alpha = 0$ ,  $\mathbf{v}_\alpha^* = \boldsymbol{\omega}$  is a solution to (3.3),
- (ii) if  $\alpha \geq \max_k \{[\mathbf{S}\boldsymbol{\omega}]_k / f_k\}$ , 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  $0 \leq \alpha \mathbf{f}^T \mathbf{v}_\alpha^* \leq \alpha \mathbf{f}^T \boldsymbol{\omega} - (\boldsymbol{\omega} - \mathbf{v}_\alpha^*)^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}_\alpha^*)$ ,
- (iv)  $\nabla D_\alpha(\mathbf{v}_\alpha^*) \geq 0$  and  $(\mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) = 0$ ,
- (v) if  $\tilde{\mathbf{v}}_\alpha^*$  is another solution to (3.3), then  $\mathbf{S}(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*) = 0$  and  $\mathbf{f}^T \tilde{\mathbf{v}}_\alpha^* = \mathbf{f}^T \mathbf{v}_\alpha^*$ ,
- (vi) if  $[\alpha \mathbf{f} - \mathbf{S}\boldsymbol{\omega}]_k \geq 0$  and  $\mathbf{S}_{k,k} > 0$  (see Remark 3.2), then  $[\mathbf{v}_\alpha^*]_k = 0$ ,
- (vii) the map  $\alpha \mapsto D(\mathbf{v}_\alpha^*)$  is increasing, and  $\alpha \mapsto \mathbf{f}^T \mathbf{v}_\alpha^*$  is decreasing.

Since  $\mathbf{v} \geq 0$ , the term  $\mathbf{f}^T \mathbf{v}$  can be interpreted as a weighted  $\ell^1$ -type regularisation (and  $\alpha$  is the regularisation parameter). For appropriate  $\mathbf{f}$  and  $\alpha$ , we can thus expect a solution  $\mathbf{v}_\alpha^*$  to (3.3) to be sparse (see, e.g., [11]). This intuition is confirmed by Lemma 3.1-(vi); indeed, when all the diagonal entries of  $\mathbf{S}$  are strictly positive,  $\text{card}(\{k | [\alpha \mathbf{f} - \mathbf{S}\boldsymbol{\omega}]_k < 0\})$  gives an upper bound on the number of strictly positive components of  $\mathbf{v}_\alpha^*$  (notice that this bound is generally not tight); remark that assertion (vi) is a direct corollary of assertion (iv) (since  $\mathbf{S}\mathbf{v}_\alpha^* \geq 0$ ).

**Remark 3.2.** Assuming  $\mathbf{S}_{k,k} = K^2(x_k, x_k) > 0$  for all  $k \in \{1, \dots, N\}$  (what we shall denote by  $\text{diag}(\mathbf{S}) > 0$ ) is equivalent to assuming  $K(x_k, x_k) > 0$  for all  $k$  (i.e.,  $\text{diag}(\mathbf{K}) > 0$ ); 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 equation (3.2) (i.e.,  $\mu$  and  $\nu$  are discrete measures supported by  $S$ ), such a point  $x_k$  may thus be removed from the sample  $S$  without inducing any modification of the operators  $T_\mu$  and  $T_\nu$ .  $\triangleleft$

**3.3. Constrained squared-kernel-discrepancy minimisation.** Instead of considering problem (3.3), we can equivalently introduce, for  $\varkappa \geq 0$  (and, in practice,  $\varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$ , see Lemma 3.2)

$$\underset{\mathbf{v}}{\text{minimise}} D(\mathbf{v}) = \frac{1}{2}(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) \text{ subject to } \mathbf{v} \geq 0 \text{ and } \mathbf{f}^T \mathbf{v} = \varkappa. \quad (3.4)$$

Notice that problem (3.4) consists in minimising a convex function on a convex compact domain, so that a solution  $\mathbf{v}_\varkappa^*$  to problem (3.4) 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 3.2.** *Let  $\mathbf{v}_\alpha^*$  be a solution to problem (3.3) with  $\alpha \geq 0$ ; then  $\mathbf{v}_\alpha^*$  is a solution to problem (3.4) with  $\varkappa = \mathbf{f}^T \mathbf{v}_\alpha^*$ . Reciprocally, assume that  $\mathbf{v}_\varkappa^*$  is a solution to problem (3.4) with  $0 < \varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$ ; then  $\mathbf{v}_\varkappa^*$  is a solution to problem (3.3) 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 (3.3) with  $\alpha \geq \max_k \{[\mathbf{S}\boldsymbol{\omega}]_k / f_k\}$ . For  $0 \leq \varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$ , the map  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  is decreasing.*

As an interesting feature, problem (3.4) can be efficiently solved thanks to a sparse descent direction QP solver (and without storing the matrix  $\mathbf{S}$ ), like for instance the vertex-exchange strategy,

see [17, Chap. 9] and Section 6.1. A sequential strategy (based on the notion of regularisation path) for solving problems (3.3) and (3.4) is also discussed in Section 5

Notice that, in view of Lemma 3.1-(iii) and Lemma 3.2, considering  $\varkappa = \rho \mathbf{f}^T \boldsymbol{\omega}$  with  $\rho \in [0, 1]$  appears as a natural parameterisation for problem (3.4).

**4. Analogy with one-class SVM.** Following for instance [18], problems (3.3) and (3.4) can be interpreted as the dual formulations of *one-class distorted SVMs* (or discrepancy-SVMs) related to the squared kernel and the initial discrete measure  $\mu$ . Soft-margin-type extensions of the one-class SVMs introduced in Sections 4.1 and 4.2 are also discussed in Appendix B.

We denote by  $\mathcal{G}$  the RKHS associated with the squared kernel  $K^2(\cdot, \cdot)$ . We introduce the function  $g_\mu \in \mathcal{G}$ , defined by  $g_\mu(x) = \int_{\mathcal{X}} K^2(t, x) d\mu(t)$ , with  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$ ; we shall refer to  $g_\mu$  as the *primal distortion term*, see also (2.1).

**4.1. One-class SVM related to the regularised problem.** We first describe the SVM related to problem (3.3). For  $g \in \mathcal{G}$ , we consider the convex minimisation problem

$$\begin{aligned} & \underset{g}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{G}}^2 + (g|g_\mu)_{\mathcal{G}} \\ & \text{subject to} && g(x_k) \geq -\alpha f_k \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (4.1)$$

The application  $g \mapsto \|g\|_{\mathcal{G}}^2$  being strictly convex, a solution to problem (4.1) is necessarily unique.

**Lemma 4.1.** *If  $\mathbf{v}_\alpha^*$  is a solution to (3.3) 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 (4.1). For all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ , we have  $g_\alpha^*(x_k) = -\alpha f_k$ .*

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  $\check{g} = g + g_\mu \in \mathcal{G}$ , problem (4.1) leads to (up to an additive constant)

$$\begin{aligned} & \underset{\check{g}}{\text{minimise}} && \frac{1}{2} \|\check{g}\|_{\mathcal{G}}^2 \\ & \text{subject to} && \check{g}(x_k) \geq g_\mu(x_k) - \alpha f_k \text{ for all } k \in \{1, \dots, N\}, \end{aligned} \quad (4.2)$$

which is an equivalent formulation for (4.1), with solution  $\check{g}_\alpha^*(x) = \sum_{k=1}^N [\mathbf{v}_\alpha^*]_k K^2(x, x_k)$ . In view of (2.1), if we denote by  $\nu_\alpha^*$  the discrete measure supported by  $\mathcal{S}$  related to a solution  $\mathbf{v}_\alpha^*$  to problem (3.3), then  $\check{g}_\alpha^* = g_{\nu_\alpha^*}$ . Problem (4.2) thus gives an interesting insight on the role of the penalisation vector  $\mathbf{f}$ , see Section 8 for a further discussion.

**Remark 4.1.** Problem (4.1) may also be considered in the framework of (3.1), i.e., for a general measure  $\mu \in \mathcal{T}(K)$ ; however, the optimisation needs in this case to be restricted to the closed linear subspace  $\mathcal{G}_{\mathcal{S}} = \text{span}\{K_{x_k}^2\}_{k=1}^N$  of  $\mathcal{G}$ , and an additive constant relative to the term  $\|K\|_{L^2(\mu \otimes \mu)}^2$  is required.  $\triangleleft$

**4.2. One-class SVM related to the constrained problem.** We now describe the SVM related to problem (3.4). For  $g \in \mathcal{G}$  and  $\gamma \in \mathbb{R}$ , we introduce the problem

$$\begin{aligned} & \underset{g, \varkappa}{\text{minimise}} && \frac{1}{2} \|g\|_{\mathcal{G}}^2 + (g|g_\mu)_{\mathcal{G}} - \varkappa \\ & \text{subject to} && g(x_k) \geq \varkappa f_k / \varkappa \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (4.3)$$

Again, a solution to problem (4.3) is necessarily unique.

**Lemma 4.2.** *If  $\mathbf{v}_\varkappa^*$  is a solution to (3.4), 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 (4.3). For all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\varkappa^*]_k > 0$ , we have  $g_\varkappa^*(x_k) = \gamma_\varkappa^* f_k / \varkappa$ .*

From Lemma 4.2, we have  $\gamma_\varkappa^* = \|g_\varkappa^*\|_{\mathcal{G}}^2 + (g_\varkappa^*|g_\mu)_{\mathcal{G}}$ . In view of Lemma 3.2, for  $0 < \varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$ , we know that  $\mathbf{v}_\varkappa^*$  is a solution to (3.3) for  $\alpha = -\gamma_\varkappa^* / \varkappa$ , and since  $\alpha \geq 0$ , we therefore have  $\gamma_\varkappa^* \leq 0$ .

**5. Regularisation path.** Considering the framework of Section 3 and following the classical results relative to the regularisation paths for Lasso or SVM models (see e.g., [15, 10]), we now discuss the regularisation paths related to problems (3.3) and (3.4). In what follows, we mainly consider problem (3.3) (i.e., the regularised model); results related to problem (3.4) can then be obtained from Lemma 3.2.



**5.1. Generalities.** Let  $\mathbf{v}_\alpha^*$  be a solution to (3.3) 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,  $[\nabla D_\alpha(\mathbf{v}_\alpha^*)]_k > 0$  for all  $k \in J_\alpha^c$ . From Lemma 3.1, the index set  $J_\alpha$  is unique, even when the solution to (3.3) 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 (3.3) for  $\alpha \geq 0$ . We also recall that for all  $k$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ , we necessarily have  $k \in J_\alpha$ , see for instance Lemma 3.1-(iv).

Knowing  $J_\alpha$ , a solution  $\mathbf{v}_\alpha^* \geq 0$  to (3.3) is characterised by the conditions

$$[\mathbf{v}_\alpha^*]_{J_\alpha^c} = 0, \text{ and } \mathbf{S}_{J_\alpha, J_\alpha} [\mathbf{v}_\alpha^*]_{J_\alpha} = [\mathbf{S}\boldsymbol{\omega}]_{J_\alpha} - \boldsymbol{\alpha} \mathbf{f}_{J_\alpha}, \quad (5.1)$$

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,  $\mathbf{f}_{J_\alpha} \in \mathbb{R}^{n_\alpha}$  stands for the vector defined by the components of  $\mathbf{f}$  with index in  $J_\alpha$ .

**Lemma 5.1.** *Let  $\mathbf{v}_{\alpha_1}^*$  and  $\mathbf{v}_{\alpha_2}^*$  be solutions to problem (3.3) with  $\alpha_1$  and  $\alpha_2 \geq 0$ , respectively. Assume that  $J_{\alpha_1} = J_{\alpha_2} = J$ , 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 (3.3) with  $\alpha = \theta \alpha_1 + (1 - \theta) \alpha_2$  and  $J_\alpha = J$ .*

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 5.1, 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 (3.3); by uniqueness of the solution, this entails that the map  $R : \alpha \mapsto \mathbf{v}_\alpha^*$  is continuous. In case of non-uniqueness of the solution, Lemma 5.1 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*, and we call *kinks* the values of  $\alpha$  where an event occurs. In particular, notice that since there cannot exist more than  $2^N$  different subsets of  $\{1, \dots, N\}$ , Lemma 5.1 implies that the number  $M_{ev}$  of events related to problem (3.3) necessarily satisfies  $M_{ev} \leq 2^N - 1$  (i.e.,  $M_{ev}$  is the number of kinks).

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 3.1, we know that for  $\alpha \geq \max_k [\mathbf{S}\boldsymbol{\omega}]_k / f_k$ , we have  $\mathbf{v}_\alpha^* = 0$ . Therefore, the first event occurs at  $\alpha_0 = \max_k [\mathbf{S}\boldsymbol{\omega}]_k / f_k$ , and we have  $J_{\alpha_0} = \{k | [\mathbf{S}\boldsymbol{\omega}]_k / f_k = \alpha_0\}$  (for larger values of  $\alpha$ , the underlying sparsity pattern is the empty set). As detailed in the following Section 5.2, 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 obtain the underlying sparsity pattern  $J_{\alpha_1}$ .

**Remark 5.1.** Lemma 5.1 can be generalised to regularised problems involving a general penalisation term  $\mathbf{r} \in \mathbb{R}^N$ , with  $\mathbf{r} \geq 0$ , i.e.,

$$\underset{\mathbf{v}}{\text{minimise}} D_{\mathbf{r}}(\mathbf{v}) = \frac{1}{2}(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) + \mathbf{r}^T \mathbf{v} \text{ subject to } \mathbf{v} \geq 0. \quad (5.2)$$

We denote by  $\mathbf{v}_{\mathbf{r}}^*$  a solution to (5.2). Note that problem (3.3) thus consists in the particular case  $\mathbf{r} = \boldsymbol{\alpha} \mathbf{f}$ . We can then define  $J_{\mathbf{r}} = \{k | [\nabla D_{\mathbf{r}}(\mathbf{v}_{\mathbf{r}}^*)]_k = 0\}$  and  $J_{\mathbf{r}}^c = \{1, \dots, N\} \setminus J_{\mathbf{r}}$ , and knowing  $J_{\mathbf{r}}$ , a solution  $\mathbf{v}_{\mathbf{r}}^* \geq 0$  to (5.2) is characterised by the conditions

$$[\mathbf{v}_{\mathbf{r}}^*]_{J_{\mathbf{r}}^c} = 0, \text{ and } \mathbf{S}_{J_{\mathbf{r}}, J_{\mathbf{r}}} [\mathbf{v}_{\mathbf{r}}^*]_{J_{\mathbf{r}}} = [\mathbf{S}\boldsymbol{\omega}]_{J_{\mathbf{r}}} - \mathbf{r}_{J_{\mathbf{r}}}$$

From exactly the same argument as in Lemma 5.1, we obtain that if two penalisation terms  $\mathbf{r}_1$  and  $\mathbf{r}_2 \geq 0$  are such that  $J_{\mathbf{r}_1} = J_{\mathbf{r}_2} = J$ , then for all  $\theta \in [0, 1]$ ,  $\mathbf{v}_{\mathbf{r}}^* = \theta \mathbf{v}_{\mathbf{r}_1}^* + (1 - \theta) \mathbf{v}_{\mathbf{r}_2}^*$  is a solution to problem (5.2) with  $\mathbf{r} = \theta \mathbf{r}_1 + (1 - \theta) \mathbf{r}_2$  and  $J_{\mathbf{r}} = J$ .  $\triangleleft$

**5.2. 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  $\alpha_0 = \max_k [\mathbf{S}\boldsymbol{\omega}]_k / f_k > \alpha_1 > \dots > \alpha_{M_{ev}-1} \geq 0$ , 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  $p \in \{0, \dots, M_{ev} - 2\}$ , i.e.,  $J_{\alpha_p} \neq \{1, \dots, N\}$ , i.e.,  $n_{\alpha_p} = \text{card}(J_{\alpha_p}) < N$ ) 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  $J_{\alpha_{p+1}}$ . 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 (5.1), we introduce the vector  $\mathbf{v}_\alpha$  such that  $[\mathbf{v}_\alpha]_{J^c} = 0$  and  $[\mathbf{v}_\alpha]_J = \mathbf{S}_{J,J}^{-1}([\mathbf{S}\boldsymbol{\omega}]_J - \boldsymbol{\alpha}\mathbf{f}_J)$ ; 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

$$[\mathbf{v}_\alpha]_J \geq 0 \text{ and } [\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \boldsymbol{\alpha}\mathbf{f}]_{J^c} \geq 0. \quad (5.3)$$

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 Lemma 5.2 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.

**Lemma 5.2.** *Consider a kink  $\alpha_p$  with related sparsity pattern  $J_{\alpha_p} = J \neq \{1, \dots, N\}$ , and assume that the submatrix  $\mathbf{S}_{J,J}$  is invertible. We introduce the  $(N - n_{\alpha_p}) \times n_{\alpha_p}$  matrix  $\mathbf{M} = \mathbf{S}_{J^c,J} \mathbf{S}_{J,J}^{-1}$ , and we define*

$$\begin{aligned} \alpha_+ &= \max_l \left\{ \left[ \mathbf{M}[\mathbf{S}\boldsymbol{\omega}]_J - [\mathbf{S}\boldsymbol{\omega}]_{J^c} \right]_l / \left[ \mathbf{M}\mathbf{f}_J - \mathbf{f}_{J^c} \right]_l \mid \left[ \mathbf{M}\mathbf{f}_J - \mathbf{f}_{J^c} \right]_l < 0 \right\}, \text{ and} \\ \alpha_- &= \max_m \left\{ \left[ \mathbf{S}_{J,J}^{-1}[\mathbf{S}\boldsymbol{\omega}]_J \right]_m / \left[ \mathbf{S}_{J,J}^{-1}\mathbf{f}_J \right]_m \mid \left[ \mathbf{S}_{J,J}^{-1}\mathbf{f}_J \right]_m < 0 \right\}. \end{aligned}$$

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.

In view of Lemma 5.2, 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  $\mathbf{S}_{J,J}^{-1}\mathbf{f}_J$  and  $\mathbf{S}_{J,J}^{-1}[\mathbf{S}\boldsymbol{\omega}]_J$  (i.e., the 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 favourable cases, the computational complexity may thus be reduced to  $\mathcal{O}(n_{\alpha_p}^2)$  (while still consisting in a direct approach). In order to further reduce the computational cost, an alternative may also consist in using an indirect iterative approach, like for example conjugate gradient methods; however, numerical errors may quickly lead to precision issues. Finally, the complexity of the two matrix vector products involving the matrix  $\mathbf{S}_{J^c,J}$  scale as  $\mathcal{O}(n_{\alpha_p}(N - n_{\alpha_p}))$ . As a result, the computation of the regularisation direction thus becomes intractable once large values of  $n_{\alpha_p}$  are reached. When  $N$  is large, the regularisation-path strategy may therefore only be used to explore the range of very sparse approximate measures  $\nu$ .

The determination of the path is also 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 9 and 10 for illustrations.

**6. Numerical solver for the constrained problem.** In this section, we discuss a numerically tractable strategy to compute approximate solutions to problem (3.4) (i.e., the constrained problem) for any given value of the parameter  $\varkappa > 0$ . We also propose two greedy exchange-type strategies aiming at enhancing the sparsity of a given approximate measure while keeping the squared-kernel discrepancy as low as possible.

**6.1. Vertex-exchange QP solver.** Consider problem (3.4); for  $\varkappa > 0$ , we can define the change of variable  $\tilde{\mathbf{v}} = \mathbf{D}\mathbf{v}$ , with  $\mathbf{D} = \text{diag}(\mathbf{d})$ , and  $\mathbf{d} = (d_1, \dots, d_N)^T = \mathbf{f}/\varkappa$ , i.e.,  $\mathbf{D}$  is a diagonal matrix with  $i$ -th diagonal entry  $d_i$  (so that  $\mathbf{d} = \mathbf{D}\mathbf{1}$ ). In this way, problem (3.4) is turned into (up to an additive constant), for  $\tilde{\mathbf{v}} \in \mathbb{R}^N$ ,

$$\text{minimise } C(\tilde{\mathbf{v}}) = \frac{1}{2}\tilde{\mathbf{v}}^T \mathbf{A}\tilde{\mathbf{v}} - \mathbf{b}^T \tilde{\mathbf{v}} \text{ subject to } \tilde{\mathbf{v}} \geq 0 \text{ and } \mathbf{1}^T \tilde{\mathbf{v}} = 1, \quad (6.1)$$

with  $\mathbf{A} = \mathbf{D}^{-1}\mathbf{S}\mathbf{D}^{-1}$  and  $\mathbf{b} = \mathbf{D}^{-1}\mathbf{S}\boldsymbol{\omega}$ . We refer to (6.1) as the *canonical QP* related to the constrained squared-kernel-discrepancy minimisation (3.4). Since  $\mathbf{A}_{i,j} = K^2(x_i, x_j)/(d_i d_j)$ , 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  $\mathbf{d}$ . Importantly, we shall therefore not store the matrix  $\mathbf{A}$ , but rather compute on the fly any required entry of  $\mathbf{A}$ ; this way, problems involving a large  $N$  may be considered. Notice that since  $\boldsymbol{\omega}$  is non-sparse, for large  $N$ , the computation of the *dual distortion term*  $\mathbf{S}\boldsymbol{\omega}$  is computationally demanding ( $\mathcal{O}(N^2)$  complexity), but it may be parallelised; approximation of the underlying kernel-matrix vector product may also be obtained by a generalized fast multipole method (see [19] and references therein). 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 polytopes defined by the constraints in (6.1) 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  $I_{\tilde{\mathbf{v}}} = \{k | \tilde{v}_k > 0\}$  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)

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

defining the sparse descent direction  $\boldsymbol{\delta} = \mathbf{e}_{i^*} - \mathbf{e}_{j^*}$  (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  $\rho$  being given by  $\rho = \min\{\tilde{v}_{j^*}, -(\boldsymbol{\delta}^T \nabla C(\tilde{\mathbf{v}}))/(\boldsymbol{\delta}^T \mathbf{A}\boldsymbol{\delta})\}$ . 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  $\nabla C(\tilde{\mathbf{v}} + \rho\boldsymbol{\delta}) = \nabla C(\tilde{\mathbf{v}}) + \rho\mathbf{A}\boldsymbol{\delta}$ , so that the gradient update involves only two columns of  $\mathbf{A}$ . The vertex-exchange strategy thus appears as a interesting candidate to scale up to relatively large  $N$  and may be used as a complement of the regularisation-path strategy described in Section 5.

Denoting by  $\tilde{\mathbf{v}}^*$  a solution to (6.1), the convergence of the vertex-exchange algorithm can be easily verified (see, e.g., [8]) 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

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

and these inequalities can also be used to check distance from optimality. In Sections 9 and 10, the accuracy of an approximate solution  $\tilde{\mathbf{v}}$  is indicated by  $\epsilon = (\tilde{\mathbf{v}} - \mathbf{e}_{i^*})^T \nabla C(\tilde{\mathbf{v}})$  (Frank-Wolfe error bound).

**6.2. Enhancing sparsity through components merging.** The canonical QP formulation introduced in Section 6.1 offers a convenient framework to enhance the sparsity of an approximate measure  $\nu$  while trying to keep its squared-kernel discrepancy as low as possible. Let  $\tilde{\mathbf{v}} \geq 0$  (with  $\tilde{\mathbf{v}} \in \mathbb{R}^N$ ) be such that  $\mathbf{1}^T \tilde{\mathbf{v}} = 1$ . In practice,  $\tilde{\mathbf{v}}$  will be an exact or approximate solution to problem (6.1), or any vector related to an interesting low-discrepancy configuration  $\mathbf{v}$  through the change of variable  $\tilde{\mathbf{v}} = \mathbf{D}\mathbf{v}$ , with  $\mathbf{D} = \text{diag}(\mathbf{d})$  and  $\mathbf{d} = \mathbf{f}/(\mathbf{f}^T \mathbf{v})$ , see Section 6.1. We assume that  $\tilde{\mathbf{v}}$  has  $n = n_0$  strictly positive components and we introduce  $I = \{i | \tilde{v}_i > 0\}$ . As illustrated in Sections 9 and 10), it is generally 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  $\tilde{\mathbf{d}}$ . We have

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

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

We can then search for the merging associated with the smallest value of  $C(\tilde{\mathbf{d}}_{\{i,j\}})$ , 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  $\{i^*, j^*\} = \operatorname{argmin}_{i \neq j} C(\tilde{\mathbf{d}}_{\{i,j\}})$ , 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  $i^* = \operatorname{argmin}_{i \neq j^*} C(\tilde{\mathbf{d}}_{\{i,j^*\}})$ , the amount of computations involved scaling as  $\mathcal{O}(n)$ .

We thus obtain a “best” pairwise merging  $\{i^*, j^*\}$  for  $\tilde{\mathbf{d}}$ . We next update all the involved objects, i.e.,  $\tilde{\mathbf{d}} \leftarrow \tilde{\mathbf{d}}_{\{i^*, j^*\}}$ ,  $I \leftarrow I \setminus \{j^*\}$ ,  $n \leftarrow n - 1$  and  $\nabla C(\tilde{\mathbf{d}}) \leftarrow \nabla C(\tilde{\mathbf{d}}_{\{i^*, j^*\}})$ , 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 tradeoff.

We thus obtain a sequence of merged vectors  $\{\tilde{\mathbf{d}}_{[0]}, \tilde{\mathbf{d}}_{[1]}, \dots, \tilde{\mathbf{d}}_{[n_0-1]}\}$ , where  $\tilde{\mathbf{d}}_0$  is our initial vector,  $\tilde{\mathbf{d}}_{[1]}$  results from the merging of two components of  $\tilde{\mathbf{d}}_{[0]}$ , etc.; by construction,  $\tilde{\mathbf{d}}_{[m]} \geq 0$  and  $\mathbb{1}^T \tilde{\mathbf{d}}_{[m]} = 1$  for all  $m$ , and  $\tilde{\mathbf{d}}_{[m]}$  has  $n_0 - m$  strictly positive components. Finally, instead of considering the approximation induced by  $\mathbf{v} = \mathbf{D}^{-1} \tilde{\mathbf{d}}_0$ , we may consider a sparser vector  $\mathbf{v}_{[m]} = \mathbf{D}^{-1} \tilde{\mathbf{d}}_{[m]}$  (notice that  $\mathbf{v}_{[m]}$  and  $\tilde{\mathbf{d}}_{[m]}$  have the same number of strictly positive components); see Sections 9 and 10 for illustrations.

**7. Approximate eigendecomposition.** This section is devoted to a discussion relative to the accuracy of the approximation of the main eigenpairs of the integral operator  $T_\mu$  induced by the eigendecomposition of the approximate operator  $T_\nu$ .

**7.1. Eigendecomposition of the approximate operator.** Let  $\nu = \sum_{k=1}^N v_k \delta_{x_k}$  be a discrete measure with support included in  $\mathcal{S}$ , with related vector  $\mathbf{v} \geq 0$ ; we assume that  $\mathbf{v} \neq \mathbf{0}$ .

We introduce the index set  $I = \{i | v_i > 0\}$  and let  $n = \operatorname{card}(I)$  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  $v_k = 0$ , since  $\nu$  is supported by only  $n$  points). Following Section 2.2, the strictly positive eigenvalues  $\{\vartheta_l | l \in \mathbb{N}_v^+\}$  of  $T_\nu$  and the associated canonically extended eigenfunctions  $\psi_l \in \mathcal{H}$ , 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 in particular obtain the o.n.b.  $\{\sqrt{\vartheta_l} \psi_l | l \in \mathbb{N}_v^+\}$  of the subspace  $\mathcal{H}_\nu$  of  $\mathcal{H}$  related to  $T_\nu$ .

As already mentioned, by definition and in view of Lemma 2.2, when  $D_{K^2}(\mu, \nu)$  is small, we can expect the main eigendirections of the operator  $T_\nu$  (the approximate eigendirection) to be relevant approximation of the main eigendirections of the operator  $T_\mu$ .

**7.2. Accuracy of the approximate eigendirections.** For simplicity, we assume that  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$ . We recall that we denote by  $\{\sqrt{\lambda_k} \varphi_k | k \in \mathbb{N}_\mu^+\}$  the o.n.b. of  $\mathcal{H}_\mu$  related to  $T_\mu$ ; the eigenfunctions  $\varphi_k$  being orthonormal in  $L^2(\mu)$ . For  $l \in \mathbb{N}_v^+$ , we introduce the *normalised approximate eigenfunctions* of  $T_\mu$  induced by  $T_\nu$ , given by

$$\hat{\varphi}_l = \psi_l / \|\psi_l\|_{L^2(\mu)}. \quad (7.1)$$

The normalised approximate eigenfunctions are such that  $\|\hat{\varphi}_l\|_{L^2(\mu)} = 1$ , but contrary to the true eigenfunctions  $\varphi_k$ , they are not necessarily orthogonal in  $L^2(\mu)$ . Notice that since  $\mathcal{H}_\nu \subset \mathcal{H}_\mu$ , we necessarily have  $\|\psi_l\|_{L^2(\mu)} > 0$ ; in addition, by definition of  $T_\mu$ ,

$$\|\psi_l\|_{L^2(\mu)}^2 = (\psi_l | T_\mu[\psi_l])_{\mathcal{H}}. \quad (7.2)$$

Controlling the orthogonality, in  $L^2(\mu)$ , between the approximate eigenfunctions  $\hat{\varphi}_l$  appears as a relatively inexpensive way to evaluate the accuracy of the approximate eigendirections. Indeed,

accurate approximate eigenfunctions  $\widehat{\varphi}_l$  should be almost mutually orthogonal in  $L^2(\mu)$ ; this condition is however only a necessary condition. It is also very instructive to try to estimate the eigenvalues, for the operator  $T_\mu$ , related to the approximate eigendirections  $\psi_l$  induced by  $T_\nu$ , as discussed hereafter.

For all  $k \in \mathbb{I}_\mu^+$ , we have  $\|\sqrt{\lambda_k} \varphi_k\|_{\mathcal{H}}^2 = 1$ . By analogy, for all  $l \in \mathbb{I}_\nu^+$ , we may define  $\widehat{\lambda}_l^{[1]}$  such that  $\|\sqrt{\widehat{\lambda}_l^{[1]}} \widehat{\varphi}_l\|_{\mathcal{H}}^2 = 1$  (so that  $\sqrt{\widehat{\lambda}_l^{[1]}} \widehat{\varphi}_l = \sqrt{\vartheta_l} \psi_l$ ). Thus, using (7.1) and (7.2), we obtain

$$\widehat{\lambda}_l^{[1]} = 1/\|\widehat{\varphi}_l\|_{\mathcal{H}}^2 = \vartheta_l \|\psi_l\|_{L^2(\mu)}^2 = (\sqrt{\vartheta_l} \psi_l | T_\mu [\sqrt{\vartheta_l} \psi_l])_{\mathcal{H}} = (T_\nu[\psi_l] | T_\mu[\psi_l])_{\mathcal{H}}. \quad (7.3)$$

If  $\psi_l$  is a true eigendirection of  $T_\mu$ , then  $\widehat{\lambda}_l^{[1]}$  corresponds to the associated eigenvalue. From the Cauchy-Schwarz inequality, we have

$$(\sqrt{\vartheta_l} \psi_l | T_\mu [\sqrt{\vartheta_l} \psi_l])_{\mathcal{H}} \leq \|\sqrt{\vartheta_l} \psi_l\|_{\mathcal{H}} \|T_\mu [\sqrt{\vartheta_l} \psi_l]\|_{\mathcal{H}} = \|T_\mu [\sqrt{\vartheta_l} \psi_l]\|_{\mathcal{H}},$$

with equality when  $\psi_l$  and  $T_\mu[\psi_l]$  are collinear, i.e., when  $\psi_l$  is a true eigendirection of  $T_\mu$ . This suggests the introduction of

$$\widehat{\lambda}_l^{[2]} = \|T_\mu [\sqrt{\vartheta_l} \psi_l]\|_{\mathcal{H}}. \quad (7.4)$$

The ratio  $\widehat{\lambda}_l^{[1]}/\widehat{\lambda}_l^{[2]}$  corresponds to the inner product, in  $\mathcal{H}$ , between the normalised vectors  $\sqrt{\vartheta_l} \psi_l$  and  $T_\mu[\sqrt{\vartheta_l} \psi_l]/\|T_\mu[\sqrt{\vartheta_l} \psi_l]\|_{\mathcal{H}}$ , and we have  $0 < \widehat{\lambda}_l^{[1]}/\widehat{\lambda}_l^{[2]} \leq 1$ . See Remark 7.1 for a discussion relative to the computation of the approximate eigenvalues  $\widehat{\lambda}_l^{[1]}$ .

Following (7.3) while considering the Hilbert structure of  $L^2(\mu)$  instead of the one of  $\mathcal{H}$ , we can also define

$$\widehat{\lambda}_l^{[3]} = (\widehat{\varphi}_l | T_\mu[\widehat{\varphi}_l])_{L^2(\mu)} = \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}^2 = (\widehat{\lambda}_l^{[2]})^2/\widehat{\lambda}_l^{[1]}, \quad (7.5)$$

so that  $\widehat{\lambda}_l^{[2]} \leq \widehat{\lambda}_l^{[3]}$ , with equality when  $\psi_l$  is an eigendirection of  $T_\mu$ . Finally, from the Cauchy-Schwarz inequality and by analogy with (7.4), we have

$$(\widehat{\varphi}_l | T_\mu[\widehat{\varphi}_l])_{L^2(\mu)} \leq \|\widehat{\varphi}_l\|_{L^2(\mu)} \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)} = \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)},$$

suggesting the introduction of

$$\widehat{\lambda}_l^{[4]} = \|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}. \quad (7.6)$$

We therefore have  $\widehat{\lambda}_l^{[3]} \leq \widehat{\lambda}_l^{[4]}$ , with, again, equality when  $\psi_l$  is an eigendirection of  $T_\mu$ . We have  $0 < \widehat{\lambda}_l^{[3]}/\widehat{\lambda}_l^{[4]} \leq 1$ , and this ratio stands for the inner product, in  $L^2(\mu)$ , between the two normalised vectors  $\widehat{\varphi}_l$  and  $T_\mu[\widehat{\varphi}_l]/\|T_\mu[\widehat{\varphi}_l]\|_{L^2(\mu)}$ .

**Remark 7.1.** For  $\mu = \sum_{k=1}^N \omega_k \delta_{x_k}$  and once  $\vartheta_l$  and  $\psi_l$  are known, the amount of computations required to compute  $\widehat{\lambda}_l^{[1]} = \vartheta_l \|\psi_l\|_{L^2(\mu)}^2$  scales as  $\mathcal{O}(N)$ , so that computing the approximate eigenvalues  $\widehat{\lambda}_l^{[1]}$  is relatively inexpensive on a numerical point of view.

Computing  $\widehat{\lambda}_l^{[2]}$ ,  $\widehat{\lambda}_l^{[3]}$  or  $\widehat{\lambda}_l^{[4]}$  requires in particular the knowledge of  $T_\mu[\psi_l]$ , which consists in performing a kernel-matrix vector product, with complexity scaling as  $\mathcal{O}(N^2)$ ; this operation is therefore costly, but may however be easily parallelised (and an approximation may for instance also be obtained by a generalised fast multipole method). Once  $T_\mu[\psi_l]$  is known, we can obtain  $\widehat{\lambda}_l^{[3]}$  and  $\widehat{\lambda}_l^{[4]}$  by computing an inner product in  $L^2(\mu)$ , with complexity scaling as  $\mathcal{O}(N)$ ; and we finally obtain  $\widehat{\lambda}_l^{[2]}$  thanks to the relation  $\widehat{\lambda}_l^{[2]} = \sqrt{\widehat{\lambda}_l^{[1]} \widehat{\lambda}_l^{[3]}}$ .

Note that computing  $\widehat{\lambda}_l^{[2]}$  directly from the Hilbert structure of  $\mathcal{H}$  requires to perform an inner product with numerical complexity scaling as  $\mathcal{O}(N^2)$ ; on a numerical point of view, this operation is therefore not interesting.  $\triangleleft$

Lemma 7.1 summarises the properties of the approximate eigenvalues  $\hat{\lambda}_l^{[1]}, \dots, \hat{\lambda}_l^{[4]}$ .

**Lemma 7.1.** *Considering equations (7.3)–(7.6), for all  $l \in \mathbb{I}_v^+$ , we have  $\hat{\lambda}_l^{[1]} \leq \hat{\lambda}_l^{[2]} \leq \hat{\lambda}_l^{[3]} \leq \hat{\lambda}_l^{[4]}$ , with equality when  $\psi_l$  is an eigendirection of  $T_\mu$ ; in case of equality, the approximation  $\hat{\lambda}_l^{[1]}$  corresponds exactly to the eigenvalue of  $T_\mu$  related to the eigendirection  $\psi_l$ .*

For  $\lambda \in \mathbb{R}$ , the function

$$\lambda \mapsto \|\lambda \sqrt{\vartheta_l} \psi_l - T_\mu[\sqrt{\vartheta_l} \psi_l]\|_{\mathcal{H}}^2 = \lambda^2 - 2\lambda \hat{\lambda}_l^{[1]} + (\hat{\lambda}_l^{[2]})^2 \quad (7.7)$$

reaches its minimum at  $\lambda = \hat{\lambda}_l^{[1]}$ ; in the same way, the function

$$\lambda \mapsto \|\lambda \hat{\varphi}_l - T_\mu[\hat{\varphi}_l]\|_{L^2(\mu)}^2 = \lambda^2 - 2\lambda \hat{\lambda}_l^{[3]} + (\hat{\lambda}_l^{[4]})^2 \quad (7.8)$$

reaches its minimum at  $\lambda = \hat{\lambda}_l^{[3]} = (\hat{\lambda}_l^{[2]})^2 / \hat{\lambda}_l^{[1]}$ .

**Remark 7.2.** Consider any measure  $\nu \in \mathcal{T}(K)$ ; the approximations  $\hat{\varphi}_l, \hat{\lambda}_l^{[1]}, \hat{\lambda}_l^{[2]}, \hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  remain unchanged if one replaces  $\nu$  by  $\theta\nu$  for any  $\theta > 0$ .  $\triangleleft$

In view of Lemma 7.1, one may assess the accuracy of the approximate eigendirections  $\psi_l$  (as eigendirections for  $T_\mu$ ) by checking how close to each other are the approximations  $\hat{\lambda}_l^{[1]}, \hat{\lambda}_l^{[2]}, \hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  (the ratios  $\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[2]}$  and  $\hat{\lambda}_l^{[3]}/\hat{\lambda}_l^{[4]}$  being of particular interest due to their geometric interpretation, as already discussed). From (7.7) and (7.8), one may for instance consider the quantities

$$[(\hat{\lambda}_l^{[2]})^2 - (\hat{\lambda}_l^{[1]})^2] / (\hat{\lambda}_l^{[1]})^2 \quad \text{and} \quad [(\hat{\lambda}_l^{[4]})^2 - (\hat{\lambda}_l^{[3]})^2] / (\hat{\lambda}_l^{[3]})^2, \quad (7.9)$$

so that the closer the (positive) terms in (7.9) are from zero, the more accurate is the approximate eigendirection  $\psi_l$ ; see Sections 9 and 10 for illustrations.

**8. Penalisation direction.** In the framework of Section 3 (i.e.,  $\nu$  is a discrete measure with support included in a fixed set of points  $\mathcal{S}$ ), sparsity of the approximate measure  $\nu$  (related to the vector  $\mathbf{v}$ ) is induced by the introduction of an  $\ell^1$ -type penalisation based on the term  $\mathbf{f}^T \mathbf{v}$  for a given  $\mathbf{f} \in \mathbb{R}^N$  with  $\mathbf{f} > 0$ . The term  $\mathbf{f}^T \mathbf{v}$  can be interpreted as the integral with respect to  $\nu$  of a (measurable) real-valued function  $f$  on  $\mathcal{X}$  satisfying  $f(x_k) = f_k$ , so that  $\mathbf{f}^T \mathbf{v} = \int_{\mathcal{X}} f(x) d\nu(x)$ . In practice, we 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, naturally raising questions relative to the choice of the penalisation direction  $\mathbf{f}$ .

**Lemma 8.1** (Penalisation direction inducing no sparsity). *If  $\mathbf{f} = \theta \mathbf{S}\boldsymbol{\omega}$  with  $\theta > 0$ , then for  $\alpha \leq 1/\theta$ ,  $\mathbf{v}_\alpha^* = (1 - \alpha\theta)\boldsymbol{\omega} \geq 0$  is a solution to (3.3) (and  $\mathbf{v}_\alpha^* = 0$  for  $\alpha > 1/\theta$ ).*

Thus, for  $\mathbf{f} = \theta \mathbf{S}\boldsymbol{\omega}$ , the solutions to (3.3) are non-sparse, and such a choice for the penalisation term  $\mathbf{f}$  is of no practical interest. More generally, we can remark that if  $\mathbf{f} = \mathbf{S}\boldsymbol{\eta} \geq 0$ , with  $\boldsymbol{\eta} \in \mathbb{R}^N$ , then for all  $\alpha$  such that  $\boldsymbol{\omega} - \alpha\boldsymbol{\eta} \geq 0$ , we have  $\nabla D_\alpha(\boldsymbol{\omega} - \alpha\boldsymbol{\eta}) = 0$ , and  $\mathbf{v}_\alpha^* = \boldsymbol{\omega} - \alpha\boldsymbol{\eta}$  is therefore a solution to (3.3).

We recall that, from (4.2), if  $\mathbf{v}_\alpha^*$  is a solution to the regularised problem (3.3) (with related measure  $\nu_\alpha^*$ ), then  $g_{\nu_\alpha^*} = \sum_{k=1}^N [\mathbf{v}_\alpha^*]_k K_{x_k}^2$  is the solution to, for  $g \in \mathcal{G}$  (with  $\mathcal{G}$  the RKHS related to the squared kernel),

$$\underset{g}{\text{minimise}} \quad \frac{1}{2} \|g\|_{\mathcal{G}}^2 \quad \text{subject to} \quad g_\mu(x_k) - g(x_k) \leq \alpha f_k \quad \text{for all } k \in \{1, \dots, N\}, \quad (8.1)$$

with  $g_\mu(x_k) = \int_{\mathcal{X}} K^2(t, x_k) d\mu(t) = [\mathbf{S}\boldsymbol{\omega}]_k$ ; in addition,  $g_\mu(x_k) - g_{\nu_\alpha^*}(x_k) = \alpha f_k$  for all  $k$  such that  $[\mathbf{v}_\alpha^*]_k > 0$ . Also, for any  $\mu$  and  $\nu \in \mathcal{T}(K)$ , we have  $D_{K^2}(\mu, \nu) = \|g_\mu - g_\nu\|_{\mathcal{G}}^2$ , and  $(g_\mu | g_\nu)_{\mathcal{G}} = \int_{\mathcal{X}} g_\mu(t) d\nu(t) = \int_{\mathcal{X}} g_\nu(t) d\mu(t)$ , see (2.1) and Appendix A.

As illustrated in Sections 9 and 10, considering  $\mathbf{f} = \mathbf{1}$  generally leads to satisfactory results; it is however also possible to consider model-dependent penalisation directions leading to interesting interpretations. Following Remark 3.2, in the framework of the regularised problem (3.3), we can reasonably assume that  $K(x_k, x_k) > 0$  for all  $k$  (i.e.,  $\text{diag}(\mathbf{K}) > 0$ ), so that, in particular,  $\mathbf{S}\boldsymbol{\omega} > 0$  (since  $\boldsymbol{\omega} > 0$ ). In what follows, we discuss penalisation schemes related to the vectors  $\mathbf{S}\boldsymbol{\omega}$  and  $\text{diag}(\mathbf{K})$ . In practice, notice that we must always ensure that the considered penalisation direction do not coincide with the pathological case described in Lemma 8.1 (for instance, if  $\mathbf{K}$  is a circulant matrix and if  $\boldsymbol{\omega} \propto \mathbf{1}$ , then  $\mathbf{S}\boldsymbol{\omega} \propto \text{diag}(\mathbf{K}) \propto \mathbf{1}$ ).

*Distortion-term-based penalisation.* In view of (8.1), considering  $\mathbf{f} = 1/(\mathbf{S}\boldsymbol{\omega})^p$  with  $p > 0$  results in a SVM model where the upper bound on  $g_\mu(x_k) - g_\nu(x_k)$  is inversely proportional to a positive power of  $g_\mu(x_k)$  (i.e., the larger  $g_\mu(x_k)$ , the smaller the bound on  $g_\mu(x_k) - g_\nu(x_k)$ ). In view of (3.3), the positive components of  $\mathbf{v}_\alpha^*$  are then more likely to correspond to large values of  $g_\mu(x_k)$ . Since, for any measure  $\nu$  supported by  $\mathcal{S}$  (with related weights  $v_k$ ), we have  $(g_\mu|g_\nu)_\mathcal{G} = \sum_{k=1}^N v_k g_\mu(x_k)$ , considering a penalisation of the form  $\mathbf{f} = 1/(\mathbf{S}\boldsymbol{\omega})^p$  with  $p > 0$  tends to promote the collinearity, in  $\mathcal{G}$ , between  $g_\mu$  and  $g_{\nu_\alpha^*}$  (by promoting a large value of their inner product).

*Kernel-diagonal-based penalisation.* We can first remark that for  $\mathbf{f} = \text{diag}(\mathbf{K})$ , we have  $\mathbf{f}^T \mathbf{v} = \text{trace}(T_\nu)$ , so that from Lemma 3.1-(iii) and by analogy with spectral truncation, a solution to the regularised problem (3.3) then satisfies  $\mathbf{f}^T \mathbf{v}_\alpha^* \leq \mathbf{f}^T \boldsymbol{\omega} = \text{trace}(T_\mu)$ .

From the reproducing property in  $\mathcal{G}$  and the Cauchy-Schwarz inequality, we have, for all  $x \in \mathcal{X}$ ,

$$\forall \mu \text{ and } \nu \in \mathcal{T}(K), |g_\mu(x) - g_\nu(x)| \leq \sqrt{D_{K^2}(\mu, \nu)} K(x, x), \quad (8.2)$$

so that for  $D_{K^2}(\mu, \nu)$  fixed,  $K(x_k, x_k)$  appears as a tight bound on  $g_\mu(x_k) - g_\nu(x_k)$ . In view of (8.1) and (8.2), by considering  $\mathbf{f} = 1/(\text{diag}(\mathbf{K}))^p$  with  $p > 0$ , we enforce the bound on the difference  $g_\mu(x_k) - g_{\nu_\alpha^*}(x_k)$  to be small at the points  $x_k$  where this difference can potentially be large, so that we can thus expect  $g_\mu(x_k) - g_{\nu_\alpha^*}(x_k)$  to be relatively small for every  $k \in \{1, \dots, N\}$ .

The impact of the penalisation direction on the tradeoff between sparsity and squared-kernel discrepancy is investigated and further discussed in the experiments of Sections 9.5 and 9.6.

**9. Two-dimensional example.** We assume that  $\mathcal{S} = \{x_k\}_{k=1}^N$  consists of the  $N = 2016$  first points of a uniform Halton sequence on  $[-1, 1]^2$  (see [14]), as illustrated in Figure 9.2. We set  $\omega_k = 1/N$  for all  $k$ , so that the measure  $\mu = \sum_k \omega_k \delta_{x_k}$  in particular 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  $\ell = 1/0.16$  (a different kernel is considered in Section 9.6). An overview of the spectrum of the operator  $T_\mu$  is given in Figure 9.1. We first consider the penalisation direction  $\mathbf{f} = \mathbf{1}$ .

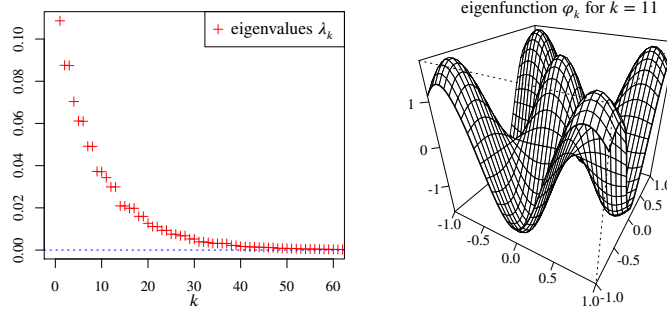


FIG. 9.1. For the two-dimensional example (Gaussian kernel and  $\boldsymbol{\omega} = \mathbf{1}/N$ ), eigenvalues  $\lambda_k$  of the integral operator  $T_\mu$  (sorted in decreasing order, only the 62 largest eigenvalues are presented), and graph, on  $[-1, 1]^2$  of the canonically extended eigenfunction  $\phi_k$  for  $k = 11$ .

**9.1. First experiment.** Figure 9.2 shows the (approximate) solution  $\mathbf{v}^*$  to problem (3.4) with  $\alpha = 0.81$ , or equivalently, to problem (3.3) with  $\alpha \approx 8.354214 \times 10^{-3}$  (with  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ). 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  $D(\mathbf{v}^*) = 7.631887 \times 10^{-4}$  (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 from the regularisation-path strategy described in Section 5 (see Section 9.2 for more details). Considering the regularisation path for problem (3.3) 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 (3.4) with increasing values of  $\kappa$ , the underlying value  $\kappa = 0.81$  satisfies

$$\kappa_p = 0.8099788 \leq \kappa \leq \kappa_{p+1} = 0.8100256, \text{ with } p = 4047.$$

In the framework of Section 6.1, the presented solution is related to a Frank-Wolfe error bound  $\epsilon = 3.989864 \times 10^{-17}$ .

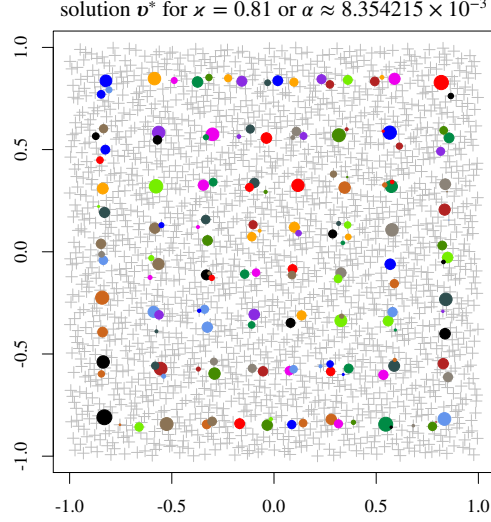


FIG. 9.2. Graphical representation (two-dimensional example, Gaussian kernel,  $\omega = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ) of the solution  $\mathbf{v}^*$  to problem (3.4) with  $\kappa = 0.81$ , or equivalently, to problem (3.3) with  $\alpha \approx 8.354215 \times 10^{-3}$ . The grey 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 accuracy of the approximate eigendecomposition of  $T_\mu$  induced by the solution  $\mathbf{v}^*$  presented in Figure 9.2 (i.e.,  $\kappa = 0.81$ ) is illustrated in Figure 9.3. In view of the similarity between the approximate eigenvalues  $\hat{\lambda}_l^{[1]}$ , and more particularly of the ratios  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[2]})^2$  and  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[3]})^2$  (see Section 7), we observe that the 21 main eigendirections of the operator  $T_{\mathbf{v}^*}$  (i.e., for  $l \in \{1, \dots, 21\}$ ) leads to remarkably accurate approximations of the eigenpairs of  $T_\mu$  related to the 21 largest eigenvalues  $\lambda_k$ . The accuracy of the approximate eigenpairs decreases for  $l \in \{22, \dots, 44\}$ , and becomes very poor for  $k \geq 44$ .

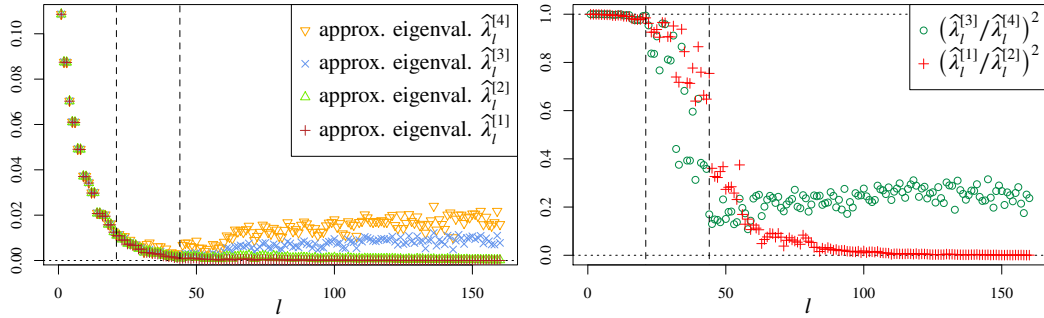


FIG. 9.3. Approximate eigenvalues  $\hat{\lambda}_l^{[1]}$ ,  $\hat{\lambda}_l^{[2]}$ ,  $\hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  induced by the solution  $\mathbf{v}^*$  presented in Figure 9.2; ratios  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[2]})^2$  and  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[3]})^2$  highlighting the accuracy of the approximate eigendirections  $\psi_l$  of  $T_\mu$  induced by  $\mathbf{v}^*$  (right).

A comparison between the true eigenvalues of  $T_\mu$  and their approximations induced by the solution  $\mathbf{v}^*$  presented in Figure 9.2 (i.e.,  $\kappa = 0.81$ ) is illustrated in Figure 9.3. We for instance observe that for  $1 \leq l \leq 8$ , the approximate eigenvalues  $\hat{\lambda}_l^{[4]}$  are the more accurate.



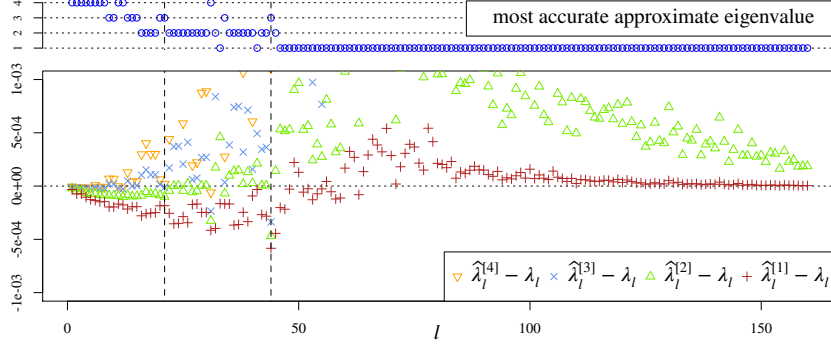


FIG. 9.4. Errors  $\hat{\lambda}_l^{[1]} - \lambda_l$  for the approximate eigenvalues induced by the solution  $\mathbf{v}^*$  presented in Figure 9.2 (bottom), and indication of the most accurate (i.e., with the smallest error in absolute value) approximate eigenvalues among  $\hat{\lambda}_l^{[1]}$ ,  $\hat{\lambda}_l^{[2]}$ ,  $\hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  (top).

**9.2. Regularisation path.** We now investigate the impact of the parameters  $\alpha$  and  $\varkappa$  related to problems (3.3) and (3.4) respectively. We compute the 12 786 first events of the regularisation-path related to problem (3.3) with decreasing values of  $\alpha$  (see Section 5), 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 (3.4) with increasing  $\varkappa$ , we have  $\varkappa_0 = 0$  and  $\varkappa_{12785} = 0.9995426$  (we recall that  $\mathbf{f}^T \boldsymbol{\omega} = 1$ ).

Figure 9.5 shows that the number of strictly positive components of the solution  $\mathbf{v}_\varkappa^*$  to problem (3.4) tends to increase when  $\varkappa$  increases. As expected from Lemma 3.1-(vii), 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 3.1 and 3.2).

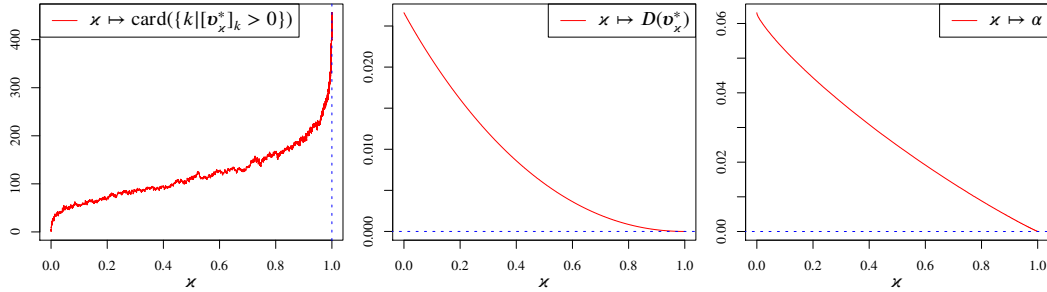


FIG. 9.5. For the two-dimensional example (Gaussian kernel,  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ), graphical representation of the 12 786 first events of the regularisation path related to problem (3.4) for increasing  $\varkappa$ ; number of strictly positive components of  $\mathbf{v}_\varkappa^*$  as function of  $\varkappa$  (left); graph of  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  (middle), and relation between  $\varkappa$  and the parameter  $\alpha$  of problem (3.3).

For 51 values of  $\varkappa$  evenly spread between 0 and  $\varkappa_{12785}$ , Figure 9.6 shows the evolution of the ratio  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[2]})^2$  for the approximate eigendecompositions induced by the 51 different solutions  $\mathbf{v}_\varkappa^*$ . As expected, the number of accurately approximate eigendirections increases with  $\varkappa$ . Remarkably, for each of the considered values of  $\varkappa$ , 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 we have  $\text{trace}(T_{\mathbf{v}_\varkappa^*}) = \varkappa$ , since  $\text{diag}(\mathbf{K}) = \mathbf{1}$  for the Gaussian kernel.

**9.3. Components merging.** We now perform the strong-pairwise-merging (see Section 6.2) of the solution  $\mathbf{v}^*$  presented in Figure 9.2 (i.e., problem (3.4) with  $\varkappa = 0.81$ ). As illustrated in Figure 9.7, 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

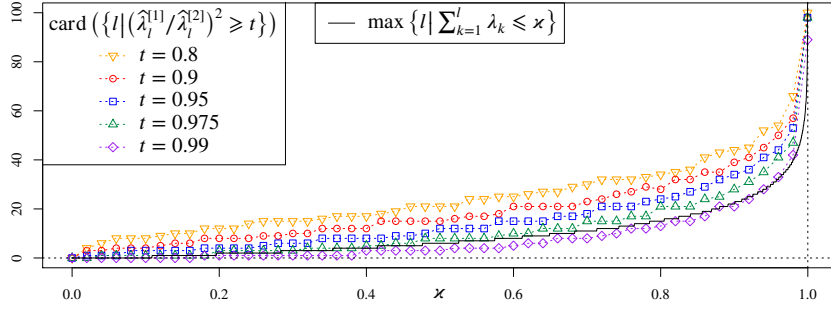


FIG. 9.6. Evolution of the accuracy of the approximate eigendecomposition of  $T_\mu$  induced by  $\mathbf{v}_x^*$  for 51 values of  $x$  between  $x_0 = 0$  and  $x_{12785} = 0.9995426$ ; the accuracy of the approximate eigendirections is measured through the ratios  $(\hat{\lambda}_l^{[1]} / \hat{\lambda}_l^{[2]})^2$ , and for illustration purpose, the map  $x \mapsto \max\{l | \sum_{k=1}^l \lambda_k \leq x\}$  is also presented (two-dimensional example, Gaussian kernel,  $\omega = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ).

the left-hand part of the figure. The accuracy of the approximate eigendecomposition induced by  $\mathbf{v}_{[90]}$  is presented in the right-hand part of Figure 9.7. 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.

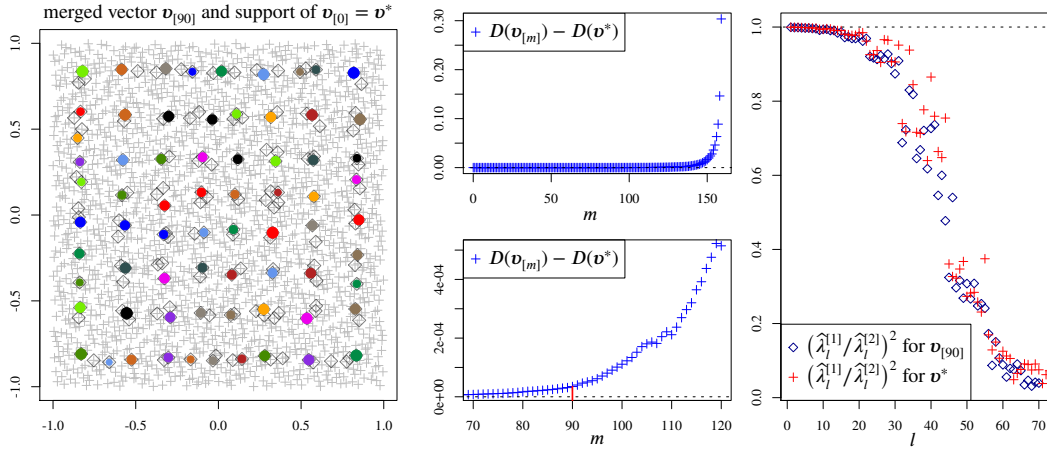


FIG. 9.7. Graphical representation of the merged solution  $\mathbf{v}_{[90]}$  (two-dimensional example with  $\omega = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ) obtained after 90 iterations of the strong-pairwise-merging strategy applied to the solution  $\mathbf{v}^*$  presented in Figure 9.2; 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 ratios  $(\hat{\lambda}_l^{[1]} / \hat{\lambda}_l^{[2]})^2$  obtained from the merged vector  $\mathbf{v}_{[90]}$  and comparison with the same ratios for the solution  $\mathbf{v}^*$  (right).

**9.4. Comparison with random sampling.** For comparison purpose, we compute the approximate eigendecompositions induced by random uniform samples (without replacement) 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 size, we perform 100 repetitions. Figure 9.8 illustrates the accuracy of the obtained approximate eigendirections, measured through the ratios  $(\hat{\lambda}_l^{[1]} / \hat{\lambda}_l^{[2]})^2$ . 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 the penalised squared-kernel-discrepancy minimisation.

**9.5. Impact of the penalisation direction.** We now study the impact of the penalisation direction  $\mathbf{f}$ . For the two-dimensional example (Gaussian kernel and  $\omega = \mathbf{1}/N$ ), we compute the

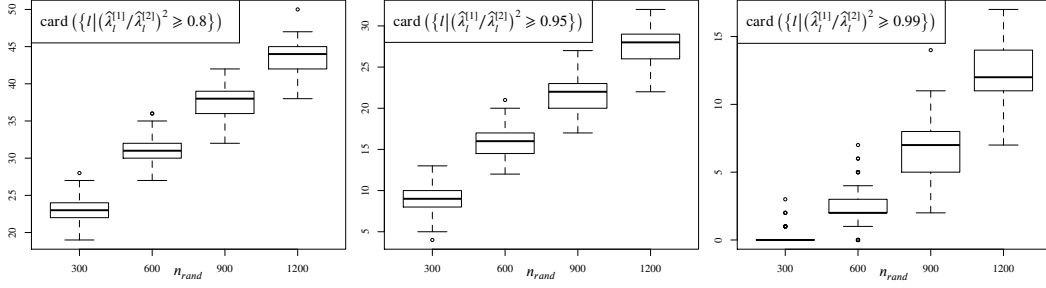


FIG. 9.8. For the two-dimensional example, accuracy of the approximate eigendecompositions induced by random samples of size  $n_{rand}$  (without replacement); for each values of  $n_{rand}$ , Tukey's boxplot, over 100 repetitions, of the number of approximate eigendirections such that  $(\hat{\lambda}_i^{[1]}/\hat{\lambda}_i^{[2]})^2 \geq 0.8$  (left), 0.95 (middle) and 0.99 (right).

regularisation path (until precision issues) of the regularised problem (3.3) for seven different vectors  $\mathbf{f} > 0$ . We consider  $\mathbf{f} = \mathbf{v}_{\max}(\mathbf{S})$  (i.e., the eigenvector related to the largest eigenvalue of the matrix  $\mathbf{S}$ , see the Perron–Frobenius theorem),  $(\mathbf{S}\boldsymbol{\omega})^2$  (i.e.,  $f_k = [\mathbf{S}\boldsymbol{\omega}]_k^2$ ),  $\sqrt{\mathbf{S}\boldsymbol{\omega}}$  (i.e.,  $f_k = \sqrt{[\mathbf{S}\boldsymbol{\omega}]_k}$ ),  $\mathbf{1}$ ,  $1/\sqrt{\mathbf{S}\boldsymbol{\omega}}$ ,  $1/(\mathbf{S}\boldsymbol{\omega})$  and  $1/(\mathbf{S}\boldsymbol{\omega})^2$ . In Figure 9.9, we represent, for each case, the number of strictly positive components of the solution to (3.3) as function of the squared-kernel discrepancy of the solution. We recall that for the Gaussian kernel, we have  $\mathbf{1} = \text{diag}(\mathbf{K}) = 1/\text{diag}(\mathbf{K})$ .

In accordance with Section 8, the results obtained for  $\mathbf{f} = \mathbf{1}$  and  $\mathbf{f} = 1/(\mathbf{S}\boldsymbol{\omega})^p$  (with in this case  $p = 1/2, 1$  and  $2$ ) appears as the more interesting in terms of the overall tradeoff between sparsity and squared-kernel discrepancy (i.e., for a same value of the squared-kernel discrepancy, the underlying solutions are generally sparser); the best results are obtained for  $\mathbf{f} = 1/(\mathbf{S}\boldsymbol{\omega})^2$ . We can however remark that  $\mathbf{f} = \mathbf{v}_{\max}(\mathbf{S})$  and  $\mathbf{f} = (\mathbf{S}\boldsymbol{\omega})^2$  lead to sparser solutions in the range of large values of the squared-kernel discrepancy (i.e., on the right-hand side of Figure 9.9).

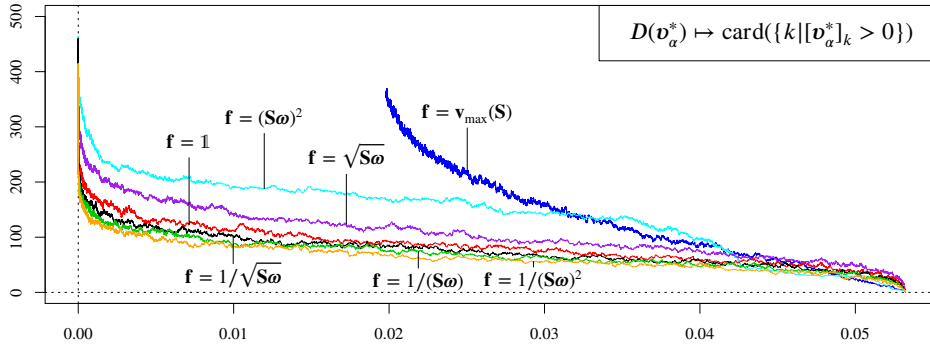


FIG. 9.9. For the two-dimensional example (Gaussian kernel and  $\boldsymbol{\omega} = \mathbf{1}/N$ ), number of strictly positive components of the solution  $\mathbf{v}_\alpha^*$  to problem (3.3) as function of the squared-kernel discrepancy  $D(\mathbf{v}_\alpha^*)$ , for various penalisation vectors  $\mathbf{f}$ ; all the curves have been obtained thanks to the regularisation path strategy (until precision issues).

**9.6. Modified kernel.** We further illustrate the impact of the penalisation direction by now considering an alternative kernel (with same set  $S$  and same initial measure  $\mu$ , i.e.,  $\boldsymbol{\omega} = \mathbf{1}/N$ , as in the previous experiments). We introduce the function, for  $x \in [-1, 1]^2$ ,  $s(x) = \sqrt{0.1 + \|x - c\|^2}$ , with  $c = (1, 1)$ , and we define the kernel (modified Gaussian kernel)

$$K(x, y) = s(x)s(y) \exp(-\ell \|x - y\|^2); \quad (9.1)$$

we still consider  $\ell = 1/0.16$ . We then in particular have  $K(x, x) = s^2(x)$ . We make the same analysis as in Section (9.5), while considering  $\mathbf{f} = \mathbf{1}$ ,  $\text{diag}(\mathbf{K})$ ,  $1/\text{diag}(\mathbf{K})$ ,  $1/(\mathbf{S}\boldsymbol{\omega})$ ,  $1/(\mathbf{S}\boldsymbol{\omega})^2$  and  $(\mathbf{S}\boldsymbol{\omega})^2$ . The results are presented in Figure 9.10. The overall tradeoff between sparsity and squared-kernel discrepancy obtained for  $\mathbf{f} = (\mathbf{S}\boldsymbol{\omega})^2$  is very poor in comparison with the tradeoffs obtained

for the five other penalisation directions, in accordance with the discussion made Section 8. The best overall tradeoffs are obtained for  $\mathbf{f} = 1/(\mathbf{S}\boldsymbol{\omega})$  and  $1/(\mathbf{S}\boldsymbol{\omega})^2$ . Among the five “efficient penalisation directions”, the trace  $\mathbf{f} = \text{diag}(\mathbf{K})$  appears to be the less interesting in the range of large values of the squared-kernel discrepancy (right-hand side of the graph), but becomes the more efficient in the range of small values of the squared-kernel discrepancy (left-hand side of the graph).

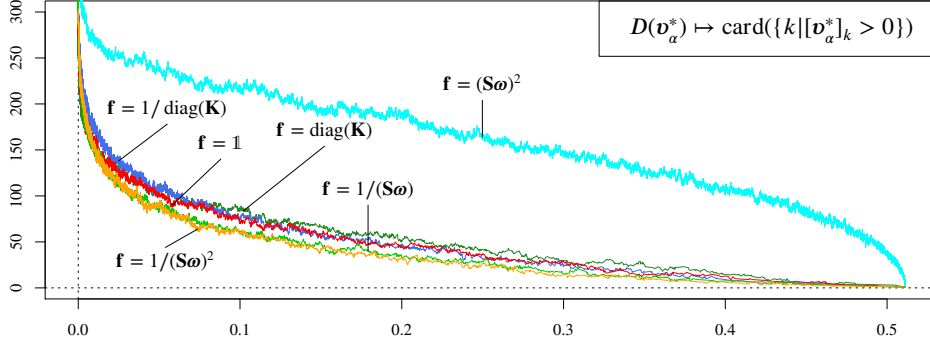


FIG. 9.10. For the two-dimensional example (modified Gaussian kernel (9.1) and  $\boldsymbol{\omega} = \mathbf{1}/N$ ), number of strictly positive components of the solution  $\mathbf{v}_\alpha^*$  to problem (3.3) as function of the squared-kernel discrepancy  $D(\mathbf{v}_\alpha^*)$ , for various penalisation vectors  $\mathbf{f}$ ; all the curves have been obtained thanks to the regularisation path strategy (until precision issues).

**10. Application to medium/large-scale problems.** This section aims at illustrating the ability of the proposed framework to tackle relatively large-scale problems. The datasets have been obtained from the UCI Machine Learning Repository, see [13]. 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 implemented in C.

**10.1. 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 9) with  $\ell = 0.02$ , and we set  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ . Notice that  $\ell = 0.02$  belongs to the range of “good parameters” for the SVM binary classification of this dataset.

We compute the 3 000 first events of the regularisation path related to problems (3.3) and (3.4). We have  $\alpha_0 = 0.2188961$  and  $\alpha_{2999} = 3.546703 \times 10^{-3}$ , and correspondingly  $\varkappa_0 = 0$  and  $\varkappa_{2999} = 0.655808$  (note that  $\mathbf{f}^T \boldsymbol{\omega} = 1$ ). A graphical representation of the obtained results is proposed in Figure 10.1. We can observe that for  $\varkappa \geq 0.5$ , the number of strictly positive components of  $\mathbf{v}_\varkappa^*$  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 took around 3 hours on our aforementioned 2015 desktop).

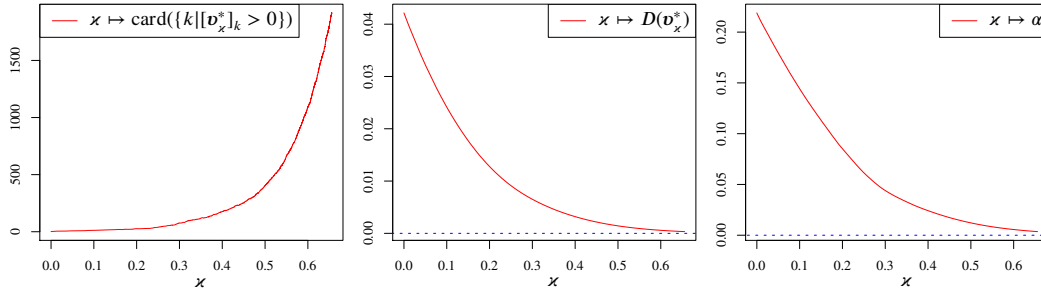


FIG. 10.1. For the MiniBooNE dataset (Gaussian kernel,  $\boldsymbol{\omega} = \mathbf{1}/N$  and  $\mathbf{f} = \mathbf{1}$ ), graphical representation of the 3 000 first events of the regularisation path related to problem (3.4) for increasing  $\varkappa$ : number of strictly positive components of  $\mathbf{v}_\varkappa^*$  as function of  $\varkappa$  (left); graph of  $\varkappa \mapsto D(\mathbf{v}_\varkappa^*)$  (middle), and relation between  $\varkappa$  and the parameter  $\alpha$  of problem (3.3).

We use the regulation path to compute the solutions to problem (3.4) for  $\varkappa = 0.3$  and  $\varkappa = 0.655$

(i.e., for problem (3.3),  $\alpha \approx 4.400276 \times 10^{-2}$  and  $\alpha \approx 3.571413 \times 10^{-3}$ ). The considered solutions have 76 and 1 902 strictly positive components, respectively. The efficiency of the induced approximate eigendecompositions is illustrated in Figure 10.2. For  $\kappa = 0.3$ , we obtain a relatively 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  $\kappa = 0.655$ , the eight main eigendirections of  $T_\mu$  are approximated with high accuracy (i.e.,  $1 \leq l \leq 8$ ), and the approximations remains relatively accurate until  $l = 29$ . Interestingly, we observe that contrary to the ratios  $(\hat{\lambda}_l^{[3]} / \hat{\lambda}_l^{[4]})^2$ , the ratios  $(\hat{\lambda}_l^{[1]} / \hat{\lambda}_l^{[2]})^2$  remain relatively high for all the values of  $l$  presented in the graph.

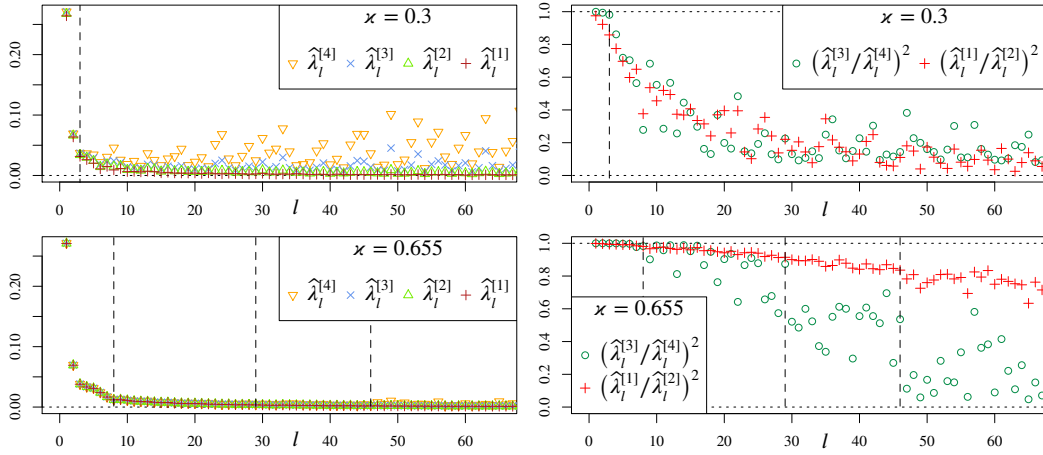


FIG. 10.2. For the MiniBooNE dataset (Gaussian kernel,  $\omega = \mathbb{1}/N$  and  $\mathbf{f} = \mathbb{1}$ ), approximate eigenvalues  $\hat{\lambda}_l^{[1]}$ ,  $\hat{\lambda}_l^{[2]}$ ,  $\hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  induced by the solution to problem (3.3) with  $\kappa = 0.3$  (top-left), and ratios  $(\hat{\lambda}_l^{[1]} / \hat{\lambda}_l^{[2]})^2$  and  $(\hat{\lambda}_l^{[3]} / \hat{\lambda}_l^{[4]})^2$  (top-right); same things for  $\kappa = 0.655$  (bottom-left) and (bottom-right).

To explore the type of solutions obtained for larger values of  $\kappa$ , we consider the vertex-exchange strategy described in Section 6.1. We compute an approximate solution for  $\kappa = 0.8$ ; 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 (3.3) verifies  $D(\hat{\mathbf{v}}^*) = 4.934072 \times 10^{-5}$  and has 9544 strictly positive components.

To enhance sparsity, we perform a weak-pairwise merging of the solution  $\hat{\mathbf{v}}^*$  for  $\kappa = 0.8$  (see Section 6.2; notice that performing a strong-pairwise merging is in this case numerically prohibitive). After 5044 iterations, the merged solution  $\mathbf{v}_{[5044]}$  is supported by 4500 points and  $D(\mathbf{v}_{[5044]}) = D(\hat{\mathbf{v}}^*) + 1.061787 \times 10^{-6}$  (i.e., increasing of 2.15%).

We next compute the approximate eigendecompositions induced by  $\hat{\mathbf{v}}^*$  and  $\mathbf{v}_{[5044]}$ ; the result are presented in Figure 10.3. In particular, in both case, the 31 main eigendirections of  $T_\mu$  are approximated with high accuracy. We observe also that for all the values of  $l$  presented on the graph, the approximation induced by  $\mathbf{v}_{[5044]}$  is equivalent, in terms of accuracy, to the approximation induced by  $\hat{\mathbf{v}}^*$ , while being related to a solution more than two times sparser.

**10.2. 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 9) with  $\ell = 0.4$ , and we set  $\omega = \mathbb{1}/N$  and  $\mathbf{f} = \mathbb{1}$ . The computation of the distortion term  $\mathbf{S}\omega$  took 5 665.6 seconds.

We compute an approximate solution (vertex-exchange strategy) for the constrained problem (3.4) with  $\kappa = 0.3$ ; 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  $\hat{\mathbf{v}}^*$  verifies  $D(\hat{\mathbf{v}}^*) = 3.931629 \times 10^{-5}$  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 10.1. We observe that a batch of 50 000 iterations of the vertex-exchange algorithm

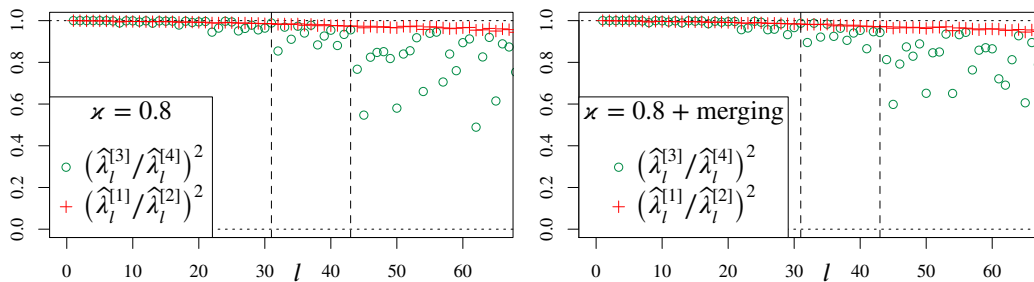


FIG. 10.3. For the MiniBooNE dataset, accuracy of the approximate eigendecompositions induced by the solution  $\hat{\mathbf{v}}^*$  to problem (3.3) with  $\kappa = 0.8$  obtained from the vertex-exchange algorithm (left), and from the merged solution  $\mathbf{v}_{[5044]}$  (left).

took around 19 minutes; the approximate solution obtained at the end of the first batch is already relatively accurate.

TABLE 10.1

For the test subsample of the SUSY dataset, information relative to the approximate solutions to problem (3.4) with  $\kappa = 0.3$  returned by the vertex-exchange algorithm for four consecutive batches of 50 000 iterations, the solver being initialised at  $\tilde{\mathbf{v}} = \mathbf{e}_1$ ; for each batch, 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                  |

To enhance the sparsity of the sampling, we perform a weak-pairwise merging of the approximate solution  $\hat{\mathbf{v}}^*$ ; the computation of 20 673 merging iterations took 78.86 seconds. The merged solution  $\mathbf{v}_{[13674]}$  is supported by 7 000 points and  $D(\mathbf{v}_{[13674]}) = D(\hat{\mathbf{v}}^*) + 5.271960 \times 10^{-7}$  (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  $\hat{\mathbf{v}}_l$  of  $\mathbf{KW}$  induced by  $\mathbf{v}_{[13674]}$  (i.e.,  $\hat{\mathbf{v}}_l \in \mathbb{R}^N$  is the vector corresponding to the approximate eigendirection  $\hat{\phi}_l$ , see Sections 2.2 and 7) took 3 313.6 seconds (time for canonical extension and rescaling), and we thus also obtain the approximate eigenvalues  $\hat{\lambda}_l^{[1]}$ . To access the accuracy of the approximate directions  $\hat{\phi}_l$ , we compute  $T_\mu[\hat{\phi}_l]$  (i.e.,  $\mathbf{KW}\hat{\mathbf{v}}_l$ ), for the 300 first approximate eigendirections, i.e.,  $l \in \{1, \dots, 300\}$ ; this operation took 191 622.3 seconds (i.e., around 53 hours). The results are presented in Figure 10.4. As already observed, the accuracy of the approximate eigendirections decreases when  $l$  increases (we recall that the eigenvalues of the approximate operator are stored in descending order). The obtained approximate eigenpairs are remarkably accurate.

**11. Conclusion.** We studied a QP-based strategy to design sparse (pointwise) quadratures for the approximation of integral operators related to symmetric positive-semidefinite kernels in a quadrature-sparsification framework (i.e., quadratures with support included in a fixed finite set of  $N$  points). For a given kernel, the considered criterion consists in the squared Hilbert-Schmidt norm for operators defined on the underlying RKHS, and sparsity of the approximate quadratures is obtained through the introduction of an  $\ell^1$ -type penalisation, under the form of a regularisation term or of a constraint. We investigated the relations between the approximation of integral operators and the notions of squared-kernel discrepancy and one-class distorted SVMs related to the squared kernel. From a methodological point of view, the considered approximation scheme is deterministic, numerically scalable (i.e., large-scale problem can be tackled) and enjoys an optimality property; it in particular



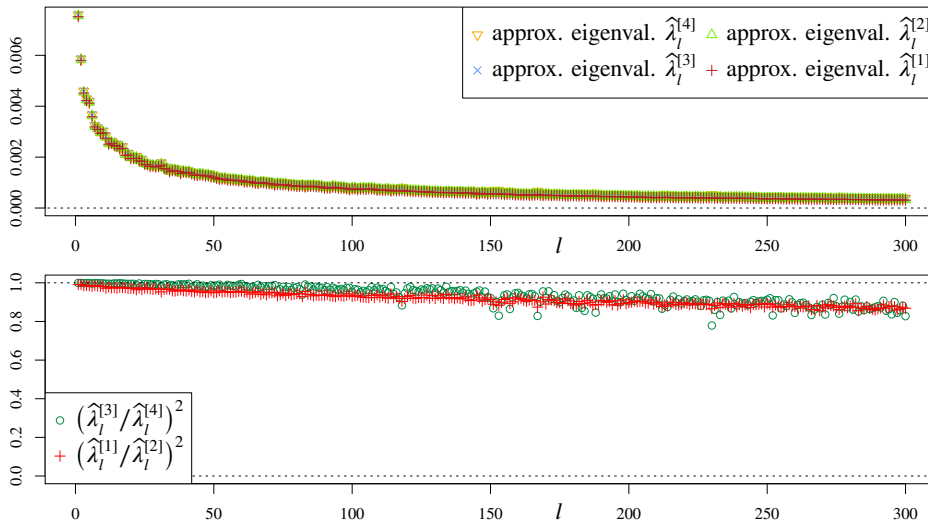


FIG. 10.4. For the test subsample of the SUSY dataset, graphical representation of the 300 first approximate eigenvalues  $\hat{\lambda}_l^{[1]}$ ,  $\hat{\lambda}_l^{[2]}$ ,  $\hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  induced by the merged solution  $\mathbf{v}_{[13674]}$  obtained from the approximate solution  $\hat{\mathbf{d}}^*$  to problem (3.4) (top); ratios  $(\hat{\lambda}_l^{[1]}/\hat{\lambda}_l^{[2]})^2$  and  $(\hat{\lambda}_l^{[3]}/\hat{\lambda}_l^{[4]})^2$  highlighting the accuracy of the underlying approximate eigendirections (bottom).

offers an interesting alternative for the sampling problem related to the Nyström approximation of large-scale kernel-matrices.

We described two numerical strategies for solving regularised or constrained squared-kernel-discrepancy minimisation 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 strategy permits the exploration of a wider range of solutions and offers a numerically efficient approach to build approximate solutions. Two greedy heuristics based on iterative pairwise-component merging are also proposed, aiming at enhancing sparsity while trying to keep squared-kernel discrepancy as low as possible. To assess the accuracy of a given approximate eigendirection, we in particular considered four different approximations (based on geometrical considerations) of the eigenvalue related to a given approximate eigendirection. We also discussed questions relative to selection of relevant penalisation directions, and we more specifically highlighted the benefits of considering penalisation terms promoting collinearity in the RKHS related to the squared kernel.

As illustrated in Section 10 (where the support of the measures  $\mu$  and  $\nu$  is included in a fixed set of  $N$  points), the proposed approach can be used to tackle relatively large-scale problems (i.e., large  $N$ ). As far as the kernel  $K(\cdot, \cdot)$  can be easily evaluated, the considered methodology leads to accurate sparse representations of the main eigenpairs of the initial operator  $T_\mu$ , in a reasonable amount of time and without necessarily resorting to powerful computing hardwares. Indeed, in the range of sparse solutions, the most computationally demanding task required to obtain an approximate measure  $\nu$  is the calculation of the distortion term  $\mathbf{S}\boldsymbol{\omega}$ , with numerical complexity scaling as  $\mathcal{O}(N^2)$ . Then, following Section 7, we can easily compute the eigendecomposition of the approximate operator  $T_\nu$ , and obtain the approximate eigenpairs  $\{(\hat{\lambda}_l^{[1]}, \hat{\varphi}_l)\}_{l \in \mathbb{I}_\nu^+}$ . Assessing the accuracy of an approximate eigendecomposition through the evaluation of the approximate eigenvalues  $\hat{\lambda}_l^{[2]}$ ,  $\hat{\lambda}_l^{[3]}$  and  $\hat{\lambda}_l^{[4]}$  is more challenging (and optional) since it requires the computation of  $T_\mu[\hat{\varphi}_l]$ , with complexity scaling as  $\mathcal{O}(N^2)$ . However, as already mentioned, any of the operations with  $\mathcal{O}(N^2)$  complexity considered in this work consists of kernel-matrix vector products that can be easily parallelised.

**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. Kernel discrepancy and integration in RKHS.

Consider the framework of Section 2 and introduce the subspace  $\mathcal{S}(K)$  of  $\mathcal{M}$ , defined as

$$\mathcal{S}(K) = \left\{ \mu \in \mathcal{M} \mid \int_{\mathcal{X}} \sqrt{K(x, x)} d\mu(x) < +\infty \right\};$$

notice that what follows may be extended to signed measures on  $\mathcal{X}$ .

From the reproducing property of the kernel  $K(\cdot, \cdot)$  and the Cauchy-Schwarz inequality, we have, for all  $h \in \mathcal{H}$  and  $\mu \in \mathcal{S}(K)$ ,

$$\left| \int_{\mathcal{X}} h(x) d\mu(x) \right| \leq \int_{\mathcal{X}} |h(x)| d\mu(x) \leq \|h\|_{\mathcal{H}} \int_{\mathcal{X}} \sqrt{K(x, x)} d\mu(x).$$

The linear functional  $I_{\mu}$  on  $\mathcal{H}$ , defined as  $I_{\mu}[h] = \int_{\mathcal{X}} h(x) d\mu(x)$ , is therefore continuous. Thus, from the Riesz representation theorem, there exists  $h_{\mu} \in \mathcal{H}$  such that  $I_{\mu}[h] = (h|h_{\mu})_{\mathcal{H}}$ , and for  $x \in \mathcal{X}$ ,  $h_{\mu}(x) = \int_{\mathcal{X}} K(x, t) d\mu(t)$ .

For  $\mu$  and  $\nu \in \mathcal{S}(K)$ , we have  $(h_{\mu}|h_{\nu})_{\mathcal{H}} = \int_{\mathcal{X} \times \mathcal{X}} K(x, t) d\mu(x) d\nu(t)$ . The *kernel discrepancy* between two measures  $\mu$  and  $\nu \in \mathcal{S}(K)$  is given by

$$D_K(\mu, \nu) = \|h_{\mu} - h_{\nu}\|_{\mathcal{H}}^2 = \|h_{\mu}\|_{\mathcal{H}}^2 + \|h_{\nu}\|_{\mathcal{H}}^2 - 2(h_{\mu}|h_{\nu})_{\mathcal{H}},$$

and  $E_K(\mu) = \|h_{\mu}\|_{\mathcal{H}}^2$  is sometime referred to as the *energy* of the measure  $\mu$  relative to  $K(\cdot, \cdot)$ .

Let us consider that we are searching for a measure  $\nu$  in order to approximate integrals relative to a fixed measure  $\mu$  (such a situation occurs for instance in case of quadrature approximation). If both measures  $\mu$  and  $\nu$  belong to  $\mathcal{S}(K)$ , from the Cauchy-Schwarz inequality, we have, for all  $h \in \mathcal{H}$ ,

$$\left| \int_{\mathcal{X}} h(x) d\mu(x) - \int_{\mathcal{X}} h(x) d\nu(x) \right| = |(h|h_{\mu} - h_{\nu})_{\mathcal{H}}| \leq \|h\|_{\mathcal{H}} \sqrt{D_K(\mu, \nu)}.$$

So, when the integrands belong to the RKHS  $\mathcal{H}$ , the error induced by approximating integrals with respect to  $\mu$  by integrals with respect to  $\nu$  has a tight bound in terms of kernel discrepancy, and it is therefore of interest to deal with a measure  $\nu$  such that  $D_K(\nu, \mu)$  is small; see [3] for a further discussion.

**Appendix B. Soft-margin-type extensions for the one-class SVM related to the regularised problem.** Pursuing the analogy with one-class SVMs, we can define soft-margin-type extensions of problems (4.1) and (4.3), 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 related to problem (4.1), but a similar discussion holds for problem (4.3).

We introduce  $\xi = (\xi_1, \dots, \xi_N)^T \in \mathbb{R}^N$ ; the components of  $\xi$  are referred to as *slack variables*. Instead of considering the constraints  $g(x_k) \geq -\alpha f_k$ , 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 [22]. 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 (3.3) (regularised squared-kernel-discrepancy minimisation), soft-margin extensions of problem (4.1) 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  $\mathbf{w}$ ).

**B.1. Weighted hinge loss.** Let  $\mathbf{c} \in \mathbb{R}^N$ , with  $\mathbf{c} \geq 0$ ; the soft-margin extension of problem (4.1) corresponding to a weighted hinge loss consists in the problem, for  $g \in \mathcal{G}$  and  $\xi \in \mathbb{R}^N$ ,

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



The Lagrange dual of problem (B.1) is given by

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

Problem (B.2) only differs from problem (3.3) 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  $\nu$ ).

**B.2. Generalised square loss.** Let  $\boldsymbol{\Sigma} \in \mathbb{R}^{N \times N}$  be a symmetric positive-definite matrix; the soft-margin extension of (4.1) corresponding to a generalised square loss is

$$\begin{aligned} & \text{minimise } \frac{1}{2} \|g\|_{\mathcal{G}}^2 + (g|g_\mu)_{\mathcal{G}} + \frac{1}{2} \boldsymbol{\xi}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\xi} \\ & \text{subject to } g(x_k) \geq -\alpha f_k - \xi_k, \text{ for all } k \in \{1, \dots, N\}. \end{aligned} \quad (\text{B.3})$$

The Lagrange dual of problem (B.3) is given by

$$\text{minimise } D_{\alpha, \boldsymbol{\Sigma}}(\mathbf{v}) = \frac{1}{2}(\boldsymbol{\omega} - \mathbf{v})^T \mathbf{S}(\boldsymbol{\omega} - \mathbf{v}) + \frac{1}{2} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} + \boldsymbol{\alpha}^T \mathbf{v} \text{ subject to } \mathbf{v} \geq 0. \quad (\text{B.4})$$

In comparison to problem (3.3), the term  $\frac{1}{2} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v}$  is added to the initial cost  $D_\alpha(\mathbf{v})$ ; it tends to “harmonise” the components of the underlying solution  $\mathbf{v}^*$ . In particular,  $D_{\alpha, \boldsymbol{\Sigma}}$  is then strongly convex.

### Appendix C. Proofs.

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

$$\begin{aligned} (T_\mu[h_j] | T_\nu[h_j])_{\mathcal{H}} &= (h_j | T_\nu[h_j])_{L^2(\mu)} = (T_\mu[h_j] | h_j)_{L^2(\nu)} \\ &= \int_{\mathcal{X} \times \mathcal{X}} K(x, t) h_j(x) h_j(t) d\mu(x) d\nu(t), \end{aligned} \quad (\text{C.1})$$

so that  $(T_\mu | T_\nu)_{\text{HS}(\mathcal{H})} = \sum_{j \in \mathbb{I}} \int_{\mathcal{X} \times \mathcal{X}} K(x, t) h_j(x) h_j(t) d\mu(x) d\nu(t)$ . For  $x$  and  $t \in \mathcal{X}$ , we have  $K(x, t) = \sum_{j \in \mathbb{I}} h_j(x) h_j(t)$ , and thus

$$\|K\|_{L^2(\mu \otimes \nu)}^2 = \int_{\mathcal{X} \times \mathcal{X}} \sum_{j \in \mathbb{I}} K(x, t) h_j(x) h_j(t) d\mu(x) d\nu(t). \quad (\text{C.2})$$

Equalities (C.1) and (C.2) hold for any o.n.b. of  $\mathcal{H}$ , so that we can in particular consider an o.n.b. which contains the o.n.b.  $\{\sqrt{\lambda_k} \varphi_k | k \in \mathbb{I}_\mu^+\}$  of  $\mathcal{H}_\mu$  defined by  $T_\mu$ . From the linear continuity of  $T_\mu$ , we then obtain

$$(T_\mu | T_\nu)_{\text{HS}(\mathcal{H})} = \sum_{k \in \mathbb{I}_\mu^+} \int_{\mathcal{X}} \lambda_k^2 \varphi_k^2(t) d\nu(t) \text{ and } \|K\|_{L^2(\mu \otimes \nu)}^2 = \int_{\mathcal{X}} \sum_{k \in \mathbb{I}_\mu^+} \lambda_k^2 \varphi_k^2(t) d\nu(t),$$

and we conclude by using the Tonelli theorem.  $\square$

*Proof of Lemma 2.2.* The proof directly follows from the properties discussed in Sections 2.1 and 2.3. In particular, (2.2) 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  $h \in \mathcal{H}_{0, \mu}$ . The inequality involving  $\tau_\mu$  is consequence of the relation  $\|h\|_{L^2(\mu)}^2 \leq \tau_\mu \|h\|_{\mathcal{H}}^2$  for all  $h \in \mathcal{H}$ .  $\square$

*Proof of Lemma 3.1.* 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 (3.3) 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 (3.3), assertion (iii) is obtained by taking  $\mathbf{v} = \boldsymbol{\omega}$ . For assertion (iv), we first remark that the first order optimality condition for  $\mathbf{v} = 0$  gives  $(\mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) \leq 0$ . If we assume that there exists  $k$  such that  $[\nabla D_\alpha(\mathbf{v}_\alpha^*)]_k < 0$ , we obtain  $(\mathbf{e}_k - \mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) < 0$ , and the first order optimality condition is then violated ( $\mathbf{e}_k$  stands for the  $k$ -th element of the canonical

basis of  $\mathbb{R}^N$ , so that  $\mathbf{e}_k^T \nabla D_\alpha(\mathbf{v}_\alpha^*) = [\nabla D_\alpha(\mathbf{v}_\alpha^*)]_k$ . Therefore, we necessarily have  $\nabla D_\alpha(\mathbf{v}_\alpha^*) \geq 0$  and  $(\mathbf{v}_\alpha^*)^T \nabla D_\alpha(\mathbf{v}_\alpha^*) = 0$  (since  $\mathbf{v}_\alpha^* \geq 0$ ). To prove (v), 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  $D_\alpha(\tilde{\mathbf{v}}_\alpha^*) = D_\alpha(\theta \mathbf{v}_\alpha^* + (1 - \theta) \tilde{\mathbf{v}}_\alpha^*)$  and  $(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \nabla D_\alpha(\tilde{\mathbf{v}}_\alpha^*) = 0$ , we necessarily have  $(\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*)^T \mathbf{S} (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*) = 0$  and therefore  $\mathbf{S} (\mathbf{v}_\alpha^* - \tilde{\mathbf{v}}_\alpha^*) = 0$  (since the matrix  $\mathbf{S}$  is symmetric and positive-semidefinite); this completes the proof. Assertion (vi) follows from the expansion, for all  $\mathbf{v} \geq 0$ ,

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

Since all the entries of  $\mathbf{S}$  are non-negative and  $\mathbf{S}_{k,k} > 0$ , if  $[\alpha \mathbf{f} - \mathbf{S} \boldsymbol{\omega}]_k \geq 0$  and  $v_k > 0$ , we necessarily have  $D(\mathbf{v} - v_k \mathbf{e}_k) < D_\alpha(\mathbf{v})$ . To obtain assertion (vii), consider  $\alpha_1 < \alpha_2$  and let  $\mathbf{v}_{\alpha_1}^*$  and  $\mathbf{v}_{\alpha_2}^*$  be solutions to problem (3.3) with  $\alpha = \alpha_1$  and  $\alpha = \alpha_2$  respectively. We have  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \leq \alpha_1 \mathbf{f}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*)$  and  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \geq \alpha_2 \mathbf{f}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*)$ , so that, necessarily,  $\mathbf{f}^T (\mathbf{v}_{\alpha_2}^* - \mathbf{v}_{\alpha_1}^*) \leq 0$ , and therefore  $D(\mathbf{v}_{\alpha_1}^*) - D(\mathbf{v}_{\alpha_2}^*) \leq 0$ .  $\square$

*Proof of Lemma 3.2.* If  $\mathbf{v}_\alpha^*$  is a solution to problem (3.3) 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  $\mathbf{f}^T \mathbf{v} = \mathbf{f}^T \mathbf{v}_\alpha^*$ , and  $\mathbf{v}_\alpha^*$  is therefore a solution to problem (3.4) with  $\varkappa = \mathbf{f}^T \mathbf{v}_\alpha^*$ .

The condition  $\varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$  follows directly from Lemma 3.1-(iii): a solution  $\mathbf{v}_\alpha^*$  to problem (3.3) indeed necessarily satisfies  $\mathbf{f}^T \mathbf{v}_\alpha^* \leq \mathbf{f}^T \boldsymbol{\omega}$ . For  $\varkappa = 0$ , we have  $\mathbf{v}_\alpha^* = 0$  and the result follows from Lemma 3.1-(ii). For  $0 < \varkappa \leq \mathbf{f}^T \boldsymbol{\omega}$ , in order to be a solution to problem (3.3), a solution  $\mathbf{v}_\alpha^*$  to problem (3.4) 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 3.1-(vii) and the correspondence between solutions to problems (3.3) and (3.4).  $\square$

*Proof of Lemma 4.1.* Define the closed linear subspace  $\mathcal{G}_S = \text{span}\{K_{x_k}^2\}_{k=1}^N$  of  $\mathcal{G}$  and let  $\mathcal{G}_0 = \mathcal{G}_S^\perp$  be its orthogonal; by definition,  $g_\mu \in \mathcal{G}_S$ . For any  $g_S \in \mathcal{G}_S$  and  $g_0 \in \mathcal{G}_0$ , we have

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

In addition, for any  $k \in \{1, \dots, N\}$ , we have  $g_0(x_k) = 0$ , so that necessarily  $g_\alpha^* \in \mathcal{G}_S$  (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 (4.1) to  $\mathcal{G}_S$  then yields, for  $\boldsymbol{\beta} \in \mathbb{R}^N$ ,

$$\text{minimise}_{\boldsymbol{\beta}} \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} \text{ subject to } \mathbf{S} \boldsymbol{\beta} \geq -\alpha \mathbf{f}. \quad (\text{C.3})$$

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

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

The primal optimality conditions give  $\mathbf{S} \boldsymbol{\beta} = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega})$ , leading to the Lagrange dual problem (3.3) (written as a minimisation problem). If  $\mathbf{v}_\alpha^*$  is a solution to (3.4), then a solution  $\boldsymbol{\beta}^*$  to (C.3) 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  $g_\alpha^* \in \mathcal{G}_S$  does not depend on such a  $\boldsymbol{\varepsilon}$ . The equality  $g_\alpha^*(x_k) = -\alpha f_k$  for all  $k \in \{1, \dots, N\}$  such that  $[\mathbf{v}_\alpha^*]_k > 0$  is consequence of the complementary slackness condition  $(\mathbf{v}_\alpha^*)^T [\mathbf{S}(\mathbf{v}_\alpha^* - \boldsymbol{\omega}) + \alpha \mathbf{f}] = 0$ .  $\square$

*Proof of Lemma 4.2.* We follow the same reasoning than in the proof of Lemma 4.1. By restricting problem (4.3) to  $\mathcal{G}_S$ , we obtain, for  $\boldsymbol{\beta} \in \mathbb{R}^N$ ,

$$\underset{\boldsymbol{\beta}, \gamma}{\text{minimise}} \quad \frac{1}{2} \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{S} \boldsymbol{\omega} - \gamma \quad \text{subject to} \quad \mathbf{S} \boldsymbol{\beta} \geq \gamma \mathbf{f} / \boldsymbol{\kappa}. \quad (\text{C.4})$$

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

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

The primal optimality conditions give  $\mathbf{S} \boldsymbol{\beta} = \mathbf{S}(\mathbf{v} - \boldsymbol{\omega})$  and  $\mathbf{f}^T \mathbf{v} = \boldsymbol{\kappa}$ , leading to the Lagrange dual problem (3.4). If  $\mathbf{v}_x^*$  is a solution to (3.4), then a solution  $\boldsymbol{\beta}^*$  to (C.4) 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  $(\mathbf{v}_x^*)^T [\mathbf{S}(\mathbf{v}_x^* - \boldsymbol{\omega}) - \gamma_x^* \mathbf{f} / \boldsymbol{\kappa}] = 0$ . The equality  $g_x^*(x_k) = \gamma_x^* f_k / \boldsymbol{\kappa}$  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 5.2.* Let  $\mathbf{v}_\alpha$  be such that  $[\mathbf{v}_\alpha]_{J^c} = 0$  and  $[\mathbf{v}_\alpha]_J = \mathbf{S}_{J,J}^{-1}([\mathbf{S}\boldsymbol{\omega}]_J - \boldsymbol{\alpha} \mathbf{f}_J)$ . Following (5.3), from the condition  $[\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \boldsymbol{\alpha} \mathbf{f}]_{J^c} \geq 0$ , we define  $\alpha_+$  as the smallest  $\alpha$  satisfying the constraint  $\alpha [\mathbf{M} \mathbf{f}_J - \mathbf{f}_{J^c}]_l \leq [\mathbf{M}[\mathbf{S}\boldsymbol{\omega}]_J - [\mathbf{S}\boldsymbol{\omega}]_{J^c}]_l$ , for all  $l \in \{1, \dots, \text{card}(J^c)\}$ . By definition, this constraint is satisfied by  $\alpha_p$ , the components  $l$  such that  $[\mathbf{M} \mathbf{d}_J - \mathbf{f}_{J^c}]_l \geq 0$  therefore carry no information, and the problem consists in searching for the smallest  $\alpha$  such that

$$\alpha \geq [\mathbf{M}[\mathbf{S}\boldsymbol{\omega}]_J - [\mathbf{S}\boldsymbol{\omega}]_{J^c}]_l / [\mathbf{M} \mathbf{f}_J - \mathbf{f}_{J^c}]_l, \quad \text{for all } l \text{ such that } [\mathbf{M} \mathbf{f}_J - \mathbf{f}_{J^c}]_l < 0.$$

In the same way, we define  $\alpha_-$  as the smallest  $\alpha$  such that  $\alpha \mathbf{S}_{J,J}^{-1} \mathbf{f}_J \leq \mathbf{S}_{J,J}^{-1} [\mathbf{S}\boldsymbol{\omega}]_J$ .  $\square$

*Proof of Lemma 5.1.* 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

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

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

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

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  $\mathbf{v}_\alpha^T (\mathbf{S}(\mathbf{v}_\alpha - \boldsymbol{\omega}) + \boldsymbol{\alpha} \mathbf{f}) = 0$ , 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 (3.3), 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. *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] Bertrand Gauthier and Luc Pronzato. Convex relaxation for IMSE optimal design in random-field models. *Computational Statistics and Data Analysis*, 2016. <http://dx.doi.org/10.1016/j.csda.2016.10.018>.
- [7] 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.
- [8] Jacques Guélat and Patrice Marcotte. Some comments on Wolfe’s “away step”. *Mathematical Programming*, 35(1):110–119, 1986.
- [9] Wolfgang Hackbusch. *Integral Equations: Theory and Numerical Treatment*, volume 120. Birkhäuser, 2012.

- [10] 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.
- [11] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical Learning with Sparsity: the Lasso and Generalizations*. CRC Press, 2015.
- [12] Sanjiv Kumar, Mehryar Mohri, and Ameet Talwalkar. Sampling methods for the Nyström method. *Journal of Machine Learning Research*, 13:981–1006, 2012.
- [13] Moshe Lichman. UCI Machine Learning Repository, 2013.
- [14] Harald Niederreiter. *Random Number Generation and Quasi-Monte Carlo Methods*. SIAM, Philadelphia, 1992.
- [15] 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.
- [16] 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.
- [17] Luc Pronzato and Andrej Pázman. *Design of Experiments in Nonlinear Models*. Springer, 2013.
- [18] 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.
- [19] Christoph Schwab and Radu Alexandru Todor. Karhunen–Loève approximation of random fields by generalized fast multipole methods. *Journal of Computational Physics*, 217(1):100–122, 2006.
- [20] Laurent Schwartz. *Analyse Hilbertienne*. Hermann, 1978.
- [21] Steve Smale and Ding-Xuan Zhou. Geometry on probability spaces. *Constructive Approximation*, 30(3):311–323, 2009.
- [22] Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer, 2008.
- [23] Shusen Wang and Zhihua Zhang. Improving CUR matrix decomposition and the Nyström approximation via adaptive sampling. *Journal of Machine Learning Research*, 14:2729–2769, 2013.