

**Sparse super resolution in microscopy: Condition,  
diffraction limit and trigonometric approximations**

**Dissertation**

ZUR ERLANGUNG DES GRADES  
DOKTOR DER NATURWISSENSCHAFTEN (DR. RER. NAT.)  
DES FACHBEREICHS MATHEMATIK/INFORMATIK/PHYSIK  
DER UNIVERSITÄT OSNABRÜCK

VORGELEGT VON

**Mathias Hockmann**

BETREUER

PROF. DR. RER. NAT. STEFAN KUNIS



OSNABRÜCK, 29. SEPTEMBER 2023



# Contents

<b>Introduction</b>	<b>5</b>
<b>1 Preliminaries</b>	<b>11</b>
1.1 Numerical linear algebra . . . . .	11
1.2 Fourier analysis . . . . .	16
1.3 Some auxiliary functions . . . . .	20
1.4 Wasserstein metric and optimal transport . . . . .	26
<b>2 Condition of sparse super resolution</b>	<b>31</b>
2.1 Summary of related results . . . . .	33
2.2 Condition estimates . . . . .	37
2.2.1 Admissible functions . . . . .	37
2.2.2 Ingham-type inequality for parameter difference . . . . .	44
2.2.3 Condition number and diffraction limit . . . . .	48
2.2.4 Condition of full inverse problem . . . . .	63
2.3 Application to Vandermonde matrices with pair clusters . . . . .	69
<b>3 Trigonometric polynomials and rational functions</b>	<b>73</b>
3.1 Approximation by polynomials . . . . .	75
3.1.1 Approximation by convolution and upper bounds . . . . .	75
3.1.2 Saturation . . . . .	79
3.1.3 Best approximation and lower bounds . . . . .	81
3.1.4 Univariate situation and uniqueness of best approximation . . . . .	84
3.2 Interpolation and approximation by the signal polynomial . . . . .	89
3.2.1 Algebraic considerations . . . . .	89
3.2.2 The signal polynomial for discrete measures . . . . .	92
3.2.3 Numerical examples for approximation and interpolation . . . . .	98
3.3 Rational Christoffel functions . . . . .	99
3.4 Recovery from noisy data . . . . .	106
<b>4 Applications in microscopy</b>	<b>115</b>
4.1 An approach to STORM analysis . . . . .	115
4.2 Structured Illumination Microscopy . . . . .	123
<b>Bibliography</b>	<b>131</b>
<b>List of figures</b>	<b>141</b>
<b>List of algorithms</b>	<b>142</b>
<b>Glossary of symbols</b>	<b>142</b>



# Introduction

The problem of recovering a highly concentrated object from blurred image data arises in many imaging applications. We abstract this as the computation of a discrete measure given access to its low-pass version or equivalently to its first Fourier coefficients and this task is then called *super resolution*. In other words, we assume access to measurements of

$$g(x) = (h * \mu)(x) = \sum_{t \in Y} \alpha_t h(x - t), \quad x \in \mathbb{R}^d \quad (1)$$

for some (de facto) bandlimited *point spread function (PSF)*  $h$  modelling the imaging device and a discrete measure  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  with support  $Y$  in some compact domain. Without loss of generality, we will consider  $Y \subset [-\frac{1}{2}, \frac{1}{2}]^2$ . While this task of recovery of  $\mu$  given  $g$  and  $h$  is already interesting for exact data from a theoretical point of view, one is confronted with noisy data in practice. Therefore, it is of great importance to control errors in the reconstruction which were caused by the noise and there has been a lot of work on the estimation of errors for specific algorithms, e.g. cf. [52, 22, 99, 101, 7, 135, 140, 49]. Beyond that, it is of equal importance to determine which amplification of the noise is inherited by the problem itself such that even the best possible algorithm cannot perform more stably than the problem allows. Even though there already exist many works in this direction of analysing the condition of the super resolution problem (see for example [53, 109, 103, 12, 37, 38, 28, 45, 44]), a major part of this dissertation is dedicated to our own approach to condition analysis overcoming several drawbacks of previous works.

In particular, this analysis allows to deepen the mathematical understanding of the various notions of a *diffraction limit* by improving a result by Chen and Moitra [28]. Historically, different diffraction limits for light microscopy were *defined* by physicists in the 19th and 20th century and we display the most well-known ones in Figure 1. There, we display the *Airy disc*, see Subsection 2.2.4 for its definition, as the prototypical PSF with some spectral bandlimit  $n$  such that *Abbe's diffraction limit* postulates that no sources with separation less than  $n^{-1}$  can be resolved (cf. [1]). In contrast to this, the Rayleigh limit [56] identifies the resolution limit at the first zero of the Airy disc leading to approximately  $1.22 \cdot n^{-1}$  as a diffraction limit while the Sparrow limit, see [147, 33], is defined as the distance such that the superposition of two translated PSFs becomes unimodal, i.e. it has just one maximum value between the two sources. Hence, this distance depends heavily on the exact knowledge of the PSF and can be numerically estimated as approximately  $0.94 \cdot n^{-1}$  for the Airy disc. Finally, the Houston criterion is an approach used in many imaging applications where the resolution is determined by the *full width at half maximum (FWHM)* of the PSFs yielding roughly  $1.03 \cdot n^{-1}$ , cf. [27, p. 47].<sup>1</sup> Even though these diffraction limits are well-established in different fields, a mathematically rigid formulation going beyond a heuristic argumentation was missing since one can argue that from a mathematical point of view two translates of  $h$  can always be distinguished from a single translate if one is given perfect measurements of  $g$  in (1). Following this reasoning, recent work of Chen and Moitra [27, 28] shows that

---

<sup>1</sup>For a more detailed description of diffraction limits see [33, 27].

the term “diffraction limit” can only be understood in the presence of noise and that this diffraction limit should then rather be seen as the point of transition from well-controllable (i.e. polynomial) to worse (i.e. exponential) noise amplification. Unfortunately, their analysis leads to a transition between  $1.15 \cdot n^{-1}$  and  $1.53 \cdot n^{-1}$  such that one cannot directly connect this to an already known and established diffraction limit.

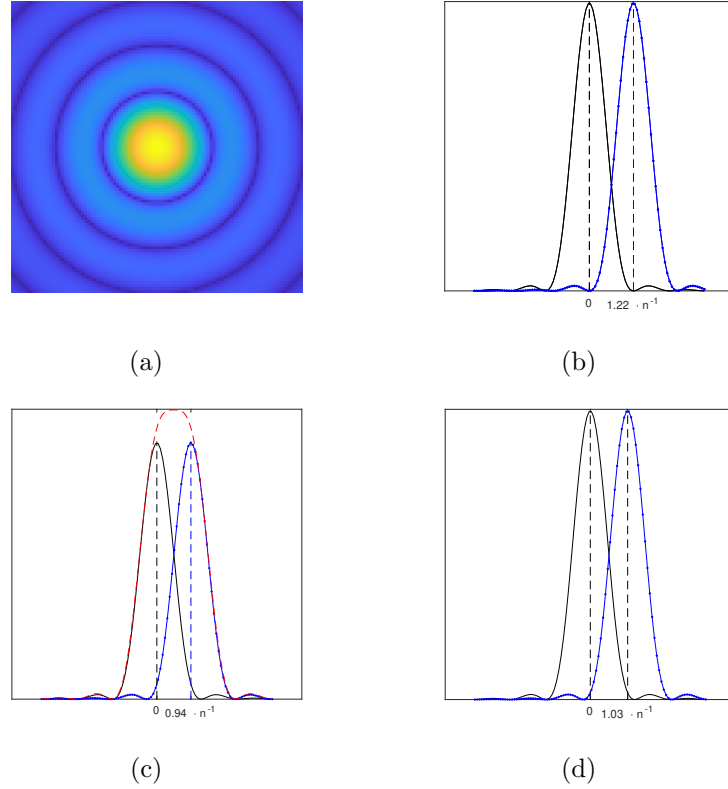


Figure 1: Various diffraction limits proposed in the literature. For a single Airy disc (a) with some bandlimit  $n$ , the Abbe diffraction limit is  $n^{-1}$  whereas the Rayleigh (b), Sparrow (c) and Houston criterion (d) provide a different reasoning yielding to different constants for the diffraction limit. For Rayleigh, two translates (black and blue) of the PSF can be separated if one is at least translated by the first zero of the PSF (b). The sparrows distance is the distance such that the superposition of the translates (red) has a single maximum instead of two local maxima (c). The Houston criterion is based on the idea that the translates are separated by the width of their main lobe measured by the *full width at half maximum* (FWHM).

Apart from being interesting in theory, bounds on the resolution are also important to guarantee the success of various algorithms solving the inverse problem of finding a discrete measure fitting to the data model (1). Usually, these methods are divided into variational methods, e.g. cf. [52, 22, 96], and subspace methods beginning with Prony [136]. More recently, machine learning approaches like [122, 121, 148] were used and implemented without any assumption on the minimal distance between support points of the measure. In [113], Mhaskar connects machine learning to the initial super resolution problem and studies the idea to approximate a measure given its low-pass data or Fourier coefficients. Because one is typically more interested in a visual representation of the support instead of

a list of molecule positions, we want to follow a similar direction and study approximations to  $\mu$  in (1) by polynomials and rational functions.

Nevertheless, each approach to solve (1) has a limited resolution by the condition of the problem itself such that there have been many attempts in microscopy to overcome the diffraction limit by changing the measurement process. Many of these approaches rely on the idea to use a non-uniform illumination of the fluorescent sample and to repeat the measurements for various illuminations. Among others, one can mention *Stochastic Optical Reconstruction Microscopy* (STORM) [139] and *Structured Illumination Microscopy* (SIM) [65, 62] in this context. While the first is considered to achieve its resolution enhancement through a repeated measurement of a random subset of the complete set of emitters, the other approach is said to improve the resolution through repeated illumination of the sample with patterned intensities. While these techniques are well-established from a practical point of view, the theoretical question of how these approaches improve the resolution is natural and its answer might allow to enhance the methods also in practice.

**Outline** After introducing main tools from linear algebra, Fourier analysis and optimal transport in Chapter 1, the subsequent chapters present the main contributions of this work. Chapter 2 deals with the analysis of the condition of sparse super resolution. After explaining connections to known results from the literature (Section 2.1), we introduce a special function in Subsection 2.2.1 that allows us to derive an inequality between the data error and the error in the parameters of the measure in Subsection 2.2.2. Such an inequality is frequently called *Ingham-inequality* in control theory, see [83, 86]. With this machinery at hand, we are able to define a condition number for the inverse problem of super resolution which leads to a diffraction limit by observing a transition in the size of the condition number from polynomial to exponential and this transition happens around the *Rayleigh criterion*  $1.22 \cdot n^{-1}$ . Interestingly, another statistical point of view leads to the same diffraction limit. After studying this in a periodic setting with Fourier coefficients of the measures as the data in Subsection 2.2.3, we show in Subsection 2.2.4 that the situation is the same for spatial data as in (1) if the PSF  $h$  satisfies some reasonable constraints. The chapter is completed by Section 2.3 where we use the Ingham inequality to obtain a new result for the smallest singular value of a Vandermonde with pairwise clustering nodes.

In Chapter 3, we study algorithmic approaches to find a measure that fits to the given Fourier data. Here, we distinguish between approximation of the measure in a weak sense (Section 3.1) and interpolation of its support by setting up a polynomial which peaks at the discrete support points of the measure and converges pointwisely to zero outside of the support as the number of moments goes to infinity, see Section 3.2. Following this idea of pointwise convergence, we study an extension to rational functions namely given by the *Christoffel function* in Section 3.3. In Section 3.4, we include the issue of noise into our assumptions on the data and show how the error in the approximations of the measure can be bounded by the size of the noise.

Finally, Chapter 4 discusses two approaches to overcome the diffraction limit found in Chapter 2. For STORM, we justify the increased resolution by a statistical argument in Section 4.1. Moreover, we apply the approaches from Chapter 3 to a large scale STORM data set and obtain a result which is comparable with typical state of the art algorithms. On the other hand, we discuss SIM and its influence to the resolution limit in Section 4.2. Starting with limitations of the popular algorithms for SIM in the case of discrete measures, we introduce an idea to circumvent these limitations by using the spatial instead of the

spectral, periodic setting where one is given shifted Fourier coefficients of the data. This allows to define a resolution limit for SIM being smaller than the previous diffraction limit from Chapter 2. Furthermore, we can rewrite the approximation algorithms from Chapter 3 for SIM and obtain their success under a weakened separation condition.

**Contributions** This dissertation contains the content of the following publications or preprints by the author in a linked and extended form:

- [25] P. Catala, M. Hockmann, and S. Kunis. Sparse super resolution and its trigonometric approximation in the p-Wasserstein distance. *Proc. Appl. Math. Mech.*, 22(1), 2023
- [26] P. Catala, M. Hockmann, S. Kunis, and M. Wageringel. Approximation and interpolation of singular measures by trigonometric polynomials. *arXiv: Numerical Analysis*, 2022
- [67] M. Hockmann and S. Kunis. Sparse super resolution is Lipschitz continuous. *arXiv: Numerical Analysis*, 2021
- [68] M. Hockmann and S. Kunis. Short Communication: Weak Sparse Superresolution is Well-Conditioned. *SIAM J. Imaging Sci.*, 16(1):SC1–SC13, 2023
- [69] M. Hockmann, S. Kunis, and R. Kurre. Towards a mathematical model for single molecule structured illumination microscopy. *Proc. Appl. Math. Mech.*, 20(1), 2021
- [70] M. Hockmann, S. Kunis, and R. Kurre. Computational resolution in single molecule localization – impact of noise level and emitter density. *Biol. Chem.*, 404(5):427–431, 2023

The main contributions of this work can be summarised as follows.

- From the technical point of view, the most important contribution is the introduction of a multivariate minorant function with specific properties, see Lemma 2.2.2, and the Ingham-type inequality in Theorem 2.2.8 being a consequence of the existence of the minorant function. Even though similar functions were known in the univariate and bivariate setting, cf. Section 1.3, our radial admissible function is not just a multivariate extension as its support in spatial domain is as small as possible, see Remark 2.2.3. Only this optimality allows to shrink the interval for the transition in the condition sufficiently in order to end up with the Rayleigh distance as a not only well-known but also mathematically natural resolution limit.
- In Theorem 2.2.21, we apply our proper definition of the condition number of the inverse problem and obtain a very sharp bounding on the diffraction limit which we define as the transition from a polynomial to an exponential condition number. This result might be seen as the most important contribution of this dissertation as it connects our condition analysis to the Rayleigh resolution limit. The emphasis on the Rayleigh limit is reinforced by an interpretation in a statistical context proving a multivariate version of an univariate conjecture from [53], see Corollary 2.2.27. As an important corollary, we provide a mathematical justification for the gain in resolution obtained by SIM (cf. Corollary 4.2.4).



- For Vandermonde matrices with pairwise clustering nodes, a lower bound on the smallest singular value is presented in Proposition 2.3.2. This result is applicable even if the cluster separation does not go to infinity as the separation within the clusters goes to zero. While there are univariate results available with this behaviour, see e.g. [11], we are not aware of any other multivariate formulation except the one presented in Proposition 2.3.2.
- Another main contribution is that studying the connection between approximation of measures and approximation of (Lipschitz) functions by polynomials leads to Wasserstein convergence rates for approximations by convolutions with kernels. In a very simple way, we can use this idea to derive an upper bound for the convergence rate in Theorem 3.1.3 whereas the lower bound in Theorem 3.1.5 and the result on the best approximation (e.g. Theorem 3.1.6) need more effort.
- We analyse how well nonlinear methods like the signal polynomial (Section 3.2) and the rational Christoffel function (Section 3.3) allow to represent a discrete measure given its moments. Main contributions in this context are the weak convergence of  $p_{1,n}$ , see Theorem 3.2.9, and the analysis of the Christoffel function for noisy data (cf. Section 3.4). In particular, we highlight the a priori choice rule for the regularisation parameter  $\varepsilon$  in Proposition 3.4.4 which guarantees that the resulting function peaks around the support of the ground truth measure  $\mu$  if the noise is sufficiently small and the separation of  $\text{supp } \mu$  is at least two times the Rayleigh condition.

**Acknowledgements** First of all, I would like to thank my supervisor Prof. Stefan Kunis for his valuable mentoring over the last years. It has always been a good mix between a relaxed working atmosphere on one hand and helpful mathematical discussions or serious “Mitarbeitergesprächen” on the other hand. During most of the time of my PhD, I obtained funding by the DFG within the Collaborative Research Center 944 “Physiology and dynamics of cellular microcompartments” for which I am very grateful. Regarding the Collaborative Research Center, I would like to thank the Group for Biophysics of the University of Osnabrück led by Prof. Jacob Piehler. Especially, I want to mention my co-supervisor Dr. Rainer Kurre who supported me in particular at the beginning of my PhD when I could always ask him questions regarding the applications in microscopy. Furthermore, I appreciate the contribution of several (anonymous) referees of my papers whose suggestions helped to improve the work presented in this dissertation. Specifically, I thank Paul Catala for proofreading of this thesis.

Even though the pandemic limited physical contacts throughout many parts of the time as a PhD student, I enjoyed the exchange with other PhD students from the institute. In the beginning, the group of the “Kaffeepause” including Ulrich, Jan-Marten, Arun, Jens, Alex, Markus, Dominik, Jonathan, Stephan and Daniel helped me at the start while the “Algebros” consisting of Anna, Harsha, Christian, Daniel, Justus, Lianne, Paul, Sarah, Tarek, Bernhard and Xiaowen made the time after the pandemic with more opportunities for social activities like the pub quiz really enjoyable. Special thanks go to Dominik, Paul and Markus for fruitful discussions about our projects, to Sarah for pointing out the connection to [18], to Tarek for enduring talks on  $q_\epsilon$  and to Daniel for the long-lasting collaboration and friendship since we met each other as freshers on our first day of studying.

Outside of university, I would like to thank my friends especially Kat, Julia, Lukas, Isabella, Yannik, Jonas, Tim, Lukas, Henning, Jan, Annelie and Lisa. All of you found

## *Introduction*

a good balance between distracting me in times where I struggled with this work or questioned myself on one hand and on the other hand asking about my work even if you were probably almost sure to get an incomprehensible answer. The same applies to my parents, Miriam, Opa Max and the rest of my family. Thanks that you are always there and that I can rely on you. Finally, I would like to express my gratitude to Merle for your encouragement over the last years. I always enjoyed to reformulate my mathematical problems into everyday language for you and you always celebrated more than I if I could solve one of those. You can find the result that we liked to call the “Bergproblem” in Lemma 2.2.4.

# 1 Preliminaries

This thesis uses methods from numerical linear algebra, Fourier analysis and optimal transport such that we spend a section for each of these areas to introduce basic notation and mention the results that are later applied in the course of this work. In between, we put a section on auxiliary functions including *minorant functions* utilised in Chapter 2.

## 1.1 Numerical linear algebra

We briefly introduce our notation for objects from linear algebra and introduce the *singular value decomposition* (SVD) as a standard tool used to extract important features of matrices. Since all presented results are well-established in the area of numerical linear algebra, there exists a wide range of literature on the topic. To give an example, we refer to [15] or [73] for an overview of (numerical) methods in linear algebra.

**Vectors, matrices and norms** We always denote vectors or functions by lower-case letters and matrices with upper-case letters. As it should be clear from the context whether object are univariate or multidimensional, we omit to use bold-face letters for vectors and matrices. The only exception is for polynomials and their coefficients where a plain letter is used for the polynomial and the same letter in bold-face for its vector of coefficients, see Section 1.3. For a complex number  $z \in \mathbb{C}$ ,  $\Re(z)$  denotes its real and  $\Im(z)$  its imaginary part. A vector  $v$  from some finite dimensional complex vector space  $\mathbb{C}^d$  is considered as a column vector where we denote by  $\bar{v}$  its complex conjugate, by  $v^\top$  its transpose and by  $v^*$  its conjugate transpose. The *inner product* between two vectors  $u, v \in \mathbb{C}^d$  with components  $(u_j)_{j=1}^d$  and  $(v_j)_{j=1}^d$  respectively is defined as

$$\langle u, v \rangle := u^* v = \sum_{j=1}^d \bar{u}_j v_j$$

leading to the *Euclidean* or *2-norm*  $\|u\|_2 = \sqrt{\langle u, u \rangle} = \left( \sum_j |u_j|^2 \right)^{1/2}$  for  $u \in \mathbb{C}^d$ . Less formally, we keep the notation shorter sporadically by denoting  $u \cdot v$  for the scalar product.<sup>2</sup> Furthermore, we use the *maximum norm* and the *p-norm* given by

$$\|u\|_\infty = \max_j |u_j| \quad \text{and} \quad \|u\|_p = \left( \sum_j |u_j|^p \right)^{1/p} \quad \text{for } p \in [1, \infty).$$

Analogously, we denote  $\bar{A}$  for the entry-wise complex conjugate of a matrix  $A \in \mathbb{C}^{m \times n}$ ,  $A^\top$  for its transpose and  $A^*$  for its conjugate transpose. The inverse of a regular square

<sup>2</sup>Additionally, we will sometimes write  $\|\cdot\|$  instead of  $\|\cdot\|_2$  if it is clear which norm is meant. Similarly, we occasionally omit the dot in the inner product and write just  $uv$  if it is clear that  $u, v$  are vectors of the same size.

## 1 Preliminaries

matrix  $A \in \mathbb{C}^{m \times m}$  is  $A^{-1}$  and  $A$  is called *Hermitian* if  $A^* = A$  or *unitary* if  $A^* = A^{-1}$ . Additionally, the eigenvalues of a square matrix  $A \in \mathbb{C}^{m \times m}$  are  $\lambda_j(A) \in \mathbb{C}$  for  $j = 1, \dots, m$  and they satisfy  $\lambda_j(A) \in \mathbb{R}$  for a Hermitian matrix  $A$ . In this case,  $\lambda_{\min}(A)$  is the smallest and  $\lambda_{\max}(A)$  is the largest eigenvalue of  $A$ . Additionally, the range or image of the matrix is denoted by  $\text{img } A$ , whereas its kernel is  $\ker A$ . The diagonal square matrix with a vector  $v = (v_j)_{j=1}^d$  on its main diagonal is represented by  $\text{diag}(v) = \text{diag}(v_1, \dots, v_d)$ . For matrices, we use the following operator norms.

**Definition 1.1.1.** (Operator norms for matrices, e.g. cf. [73, Sec. 5.6]) For  $m, n \in \mathbb{N}$  and  $A \in \mathbb{C}^{m \times n}$ , one defines<sup>3</sup>

(i) the *spectral* or *2-norm*

$$\|A\|_2 := \max_{x \in \mathbb{C}^n} \frac{\|Ax\|_2}{\|x\|_2} = \max_{x \in \mathbb{C}^n: \|x\|_2=1} \|Ax\|_2 = \sqrt{\lambda_{\max}(A^*A)},$$

(ii) the  $\infty$ -norm being the maximum  $\ell^1$ -norm of a row

$$\|A\|_\infty := \max_{x \in \mathbb{C}^n} \frac{\|Ax\|_\infty}{\|x\|_\infty} = \max_{x \in \mathbb{C}^n: \|x\|_\infty=1} \|Ax\|_\infty = \max_i \sum_j |A_{ij}|,$$

(iii) and the *1-norm* being the maximum  $\ell^1$ -norm of a column

$$\|A\|_1 := \max_{x \in \mathbb{C}^n} \frac{\|Ax\|_1}{\|x\|_1} = \max_{x \in \mathbb{C}^n: \|x\|_1=1} \|Ax\|_1 = \max_j \sum_i |A_{ij}|.$$

In contrast to this, the Froebenius norm

$$\|A\|_F = \left( \sum_{i,j} |A_{i,j}|^2 \right)^{1/2}$$

satisfies  $\|I\|_F \neq 1$  for a non-scalar identity matrix  $I$  and hence it is a matrix norm not originating from an operator norm (e.g. cf. [73]). As all norms on finite spaces like  $\mathbb{C}^d$  or  $\mathbb{C}^{m \times n}$  are equivalent, one can bound every norm for  $m \times n$ -matrices from above and below by any other norm on the space of  $m \times n$ -matrices. For example, we have the following:

**Lemma 1.1.2.** (cf. [58, Sec. 2.3]) For  $A \in \mathbb{C}^{m \times n}$  we have

(i)  $\|A\|_2 \leq \|A\|_F \leq \sqrt{\text{rank}(A)} \|A\|_2$  where  $\text{rank}(A)$  is the rank of  $A$ ,

(ii)  $m^{-1/2} \|A\|_1 \leq \|A\|_2 \leq n^{1/2} \|A\|_1$ ,

(iii)  $n^{-1/2} \|A\|_\infty \leq \|A\|_2 \leq m^{1/2} \|A\|_\infty$  and

(iv)  $\|A\|_2 \leq \sqrt{\|A\|_1 \cdot \|A\|_\infty}$ .

Matrix norms allow to study the *condition* of linear problems. More general, the condition of a possibly nonlinear map  $\phi : X \rightarrow Y$  between normed spaces  $X, Y$  measures how severely errors in the input  $x \in X$  are amplified in the output  $y = \phi(x)$ . In this sense, *conditioning* describes the behaviour of a mathematical problem encoded by  $\phi$  in

<sup>3</sup>The last equalities in (i)-(iii) can be derived from the definitions, see [73, Sec. 5.6].

the presence of noise on the data. In contrast to that, the term *stability* characterises the “perturbation behaviour of an algorithm to solve that problem on a computer” ([149, p.89]). Hence, the condition of a problem is only a property of the problem itself whereas the stability of an algorithm for the solution of the problem might be even worse. A quantification of the condition is given by a *condition number* and the literature distinguishes between the *absolute condition number* and the *relative condition number*.

**Definition 1.1.3.** (Condition number, e.g. cf. [149, p. 90]) Let  $X, Y$  be Banach spaces with a map  $\phi : X \rightarrow Y$ . Then, the absolute condition number of  $\phi$  at  $x \in X$  is

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{x': \|x-x'\| \leq \delta} \frac{\|\phi(x) - \phi(x')\|}{\|x - x'\|}$$

whereas the relative condition number considers the ratio of the relative differences in input and output space by defining

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{x': \|x-x'\| \leq \delta} \frac{\|\phi(x) - \phi(x')\| \|x\|}{\|x - x'\| \|\phi(x)\|}.$$

Often, the limit process that defines the two types of condition numbers is approximated by choosing a small  $\delta > 0$  (cf. [19, p.xix]). Whether to study the absolute or relative condition number depends on the problem. Moreover, the nature of the problem itself might lead to perturbations which have a certain structure and incorporating this leads to a *structured condition number* (cf. [19, p.119]).

Note that the definition of the absolute condition number does also make sense in metric spaces which will be the reason to study an absolute condition number in Chapter 2. If  $\phi$  is differentiable, it is easy to see that  $\kappa_{\text{abs}} = \|\phi'(x)\|$  and  $\kappa_{\text{rel}} = \frac{\|\phi'(x)\| \|x\|}{\|\phi(x)\|}$  (cf. [149, p. 90]).

**Example 1.1.4.** A fundamental condition number in numerical linear algebra arises from considering the mapping  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m, x \mapsto Ax$ , for some  $A \in \mathbb{R}^{m \times n}$ . Here we have

$$\kappa_{\text{abs}} = \lim_{\delta \rightarrow 0} \sup_{x': \|x-x'\| \leq \delta} \frac{\|A(x - x')\|}{\|x - x'\|} \leq \|A\|$$

and for  $n = m$  also

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{x': \|x-x'\| \leq \delta} \frac{\|A(x - x')\| \|A^{-1}Ax\|}{\|x - x'\| \|Ax\|} \leq \|A\| \|A^{-1}\|$$

if  $A$  is regular. For an operator norm  $\|\cdot\|$  the vector  $x$  can be chosen such that equality holds and hence the expression  $\text{cond}(A) := \|A\| \|A^{-1}\|$  is called the *condition number of the matrix*  $A$  (e.g. [149, pp. 93-94]). This definition can be extended to continuous, bijective operators between Banach spaces ([5, Thm. 2.4.3]). For a rectangular, full rank matrix  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ , the condition number is defined through the *Moore-Penrose pseudo inverse* which is part of the next paragraph.

**Singular value decomposition and its perturbation** The concept of representing a Hermitian, positive definite matrix through its eigenvalue decomposition can be generalised through the *singular value decomposition*.

**Definition 1.1.5.** (Singular value decomposition, e.g. cf. [15, Thm. 1.1.6]) A *singular value decomposition (SVD)* of a matrix  $A \in \mathbb{C}^{m \times n}$  is

$$A = U\Sigma V^* = \begin{pmatrix} U_1 & U_2 \end{pmatrix} \begin{pmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} V_1 & V_2 \end{pmatrix}^*$$

with a diagonal matrix  $\Sigma_1 = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$  where  $r = \text{rank}(A)$  and the monotony  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  holds. The matrices  $U \in \mathbb{C}^{m \times m}, V \in \mathbb{C}^{n \times n}$  are unitary. While the nonnegative values  $\sigma_j, j = 1, \dots, r$ , are called *singular values*, the columns of  $U, V$  are called *left* or *right singular vectors* respectively. Moreover, the *compact* or *truncated SVD*  $A = U_1 \Sigma_1 V_1^*$  is obtained by just using the positive singular values.

One can easily observe a connection of the SVD and the eigenvalue decomposition of  $A^*A$  by considering

$$A^*A = V\Sigma^*\Sigma V^*$$

if  $A = U\Sigma V^*$ . Therefore, a SVD of any matrix always exists and the right singular vectors are eigenvectors of  $A^*A$  while the singular values are the positive roots of eigenvalues of  $A^*A$ . Moreover, we see due to the properties of eigendecompositions of Hermitian matrices that the singular values are always unique and that the singular vectors are unique up to scalar multiplication by  $\lambda \in \mathbb{C}, |\lambda| = 1$  if they correspond to a simple singular value ([15, pp. 33-34]).

Similarly to the Courant-Fischer theorem for eigenvalues of Hermitian matrices, the variational formulation of the singular values ranks the contribution of each of the singular vectors in the representation of the matrix. We omit to present the variational formulation of every singular value and focus on the smallest and largest singular value which we often denote by  $\sigma_{\min}$  and  $\sigma_{\max}$  rather than  $\sigma_r$  and  $\sigma_1$  in order to prevent confusion whether  $\sigma_r$  or  $\sigma_1$  is larger.

**Lemma 1.1.6.** (e.g. cf. [73, Thm. 7.3.8]) For  $A \in \mathbb{C}^{m \times n}$  with  $\text{rank}(A) = n \leq m$  we have

$$\sigma_{\max}(A) = \max_{x \in \mathbb{C}^n} \frac{\|Ax\|_2}{\|x\|_2} = \|A\|_2 \quad \text{and} \quad \sigma_{\min}(A) = \min_{x \in \mathbb{C}^n} \frac{\|Ax\|_2}{\|x\|_2}$$

where the characterisation of  $\sigma_{\max}$  remains valid without the assumption on the rank of the matrix  $A$ .<sup>4</sup>

*Proof.* The proof follows directly by the min-max characterisation of singular vectors presented in [73, Thm. 7.3.8]. More precisely, studying [73, Eq. 7.3.9] for  $\sigma_1$  gives directly the result for  $\sigma_{\max}$  while the characterisation of  $\sigma_{\min}$  is an immediate consequence of [73, Eq. 7.3.10].  $\square$

**Remark 1.1.7.** Among many others, one can use the SVD for the following:

- (i) If  $A = U\Sigma V^* \in \mathbb{C}^{m \times n}$  one defines the *Moore-Penrose pseudo inverse* by

$$A^\dagger = V\Sigma^{-1}U^* \in \mathbb{C}^{n \times m} \quad \text{where} \quad \Sigma^{-1} := \begin{pmatrix} \Sigma_1^{-1} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{C}^{n \times m}$$

---

<sup>4</sup>As it should be clear from the context, we usually do not indicate explicitly from which matrix the singular values stem from and write  $\sigma_j$  instead of  $\sigma_j(A)$ .

results of inverting all (non-zero) singular values. While it specialises to the usual matrix inverse for regular square matrices, the pseudo inverse arises naturally in least squares problems since  $x^\dagger := A^\dagger b$  is the minimal norm solution of a least squares problem, i.e.  $\|x^\dagger\|_2 = \min_{x \in \arg\min \|Ax-b\|_2} \|x\|_2$  (cf. [15, Thm. 2.1.2]). Moreover, proceeding with Example 1.1.4 we can calculate for the case of the spectral norm

$$\kappa_{\text{rel}} = \lim_{\delta \rightarrow 0} \sup_{x': \|x-x'\|_2 \leq \delta} \frac{\|A(x-x')\|_2 \|x\|_2}{\|x-x'\|_2 \|Ax\|_2} \leq \|A\|_2 \frac{1}{\min \frac{\|Ax\|_2}{\|x\|_2}} = \frac{\sigma_{\max}}{\sigma_{\min}}$$

and thus the *spectral condition number* of  $A$ ,  $\text{cond}_2(A) := \|A\|_2 \cdot \|A^\dagger\|_2 = \frac{\sigma_{\max}}{\sigma_r}$ , can be defined for any  $A$  with  $\text{rank}(A) = r > 0$ .

- (ii) As the SVD of a matrix  $A$  contains information about the contribution of the subspaces spanned by the singular vectors of the matrix  $A$ , it can also be used for low rank approximation of  $A$  where one deals with the problem

$$\min_{B \in \mathbb{C}^{m \times n}, \text{rank}(B)=r_0} \|A - B\|, \quad r_0 < \min(m, n),$$

for a given matrix  $A \in \mathbb{C}^{m \times n}$  and a unitarily invariant norm  $\|\cdot\|$ , i.e.  $\|PAQ\| = \|A\|$  for any unitary matrices  $P, Q$ . According to the well-known *Eckhart-Young-Mirsky theorem* (e.g. cf. [15, Thm. 2.2.11]), an optimal solution for the spectral or the Froebenius norm is  $B = A_{r_0} = \sum_{j=1}^{r_0} \sigma_j u_j v_j^*$  where  $u_j, v_j$  are the  $j$ -columns of  $U$  and  $V$  respectively. Moreover, the residual can be expressed in terms of the singular values by

$$\|A - A_{r_0}\|_2 = \sigma_{r_0+1} \quad \text{and} \quad \|A - A_{r_0}\|_F = \sqrt{\sum_{j=r_0+1}^{\text{rank}(A)} \sigma_j^2}.$$

Even though the pseudo inverse of a product is in general not equal to product of the pseudo inverses, the following lemma about the pseudo inverse of a product with equal rank gives a sufficient condition and it is used in Chapter 3.

**Lemma 1.1.8.** (cf. [15, p. 230]) *Let  $A \in \mathbb{C}^{N \times r}$  be a matrix of rank  $r$  and let  $W \in \mathbb{C}^{r \times r}$  be a non-singular matrix. If  $T = AW A^*$ , then*

$$T^\dagger = (A^*)^\dagger W^{-1} A^\dagger.$$

*Proof.* The result can be obtained by applying [15, Thm. 2.2.3] twice in the computation

$$T^\dagger = (AW A^*)^\dagger = (W A^*)^\dagger A^\dagger = (A^*)^\dagger W^{-1} A^\dagger$$

because all matrices in the factorisation have equal rank. □

A classical result which can be used for perturbation analysis is the following.

**Lemma 1.1.9.** (Weyl's inequality cf. [15, Thm. 2.2.10]) *If  $A = B + C \in \mathbb{C}^{m \times n}$ , one has*

$$\sigma_{i+j-1}(A) \leq \sigma_i(B) + \sigma_j(C), \quad 1 \leq i + j - 1, i, j \leq n \leq m$$

*for the ordered singular values. Taking  $j = 1$  yields*

$$|\sigma_i(A) - \sigma_i(B)| \leq \sigma_1(A - B) = \|A - B\|_2.$$

## 1 Preliminaries

We remark that there are also results for the perturbation of singular vectors or more precisely for the subspaces spanned by them, e.g. cf. [155]. However, we circumvent the introduction of these subspaces by the following perturbation theorem for the pseudo inverse.

**Theorem 1.1.10.** (Wedin, cf. [156]) *If  $\text{rank } A = \text{rank } B$ , we have*

$$\|A^\dagger - B^\dagger\|_2 \leq \frac{1 + \sqrt{5}}{2} \|A^\dagger\|_2 \|B^\dagger\|_2 \|A - B\|_2.$$

## 1.2 Fourier analysis

The theory of Fourier analysis provides many tools for signal and imaging applications. Because of that and its long history, there exists a large variety of introductory literature including [61, 80, 132]. We briefly summarise some Fourier analytic results and introduce our notation for them.

**Fourier transform and function spaces on the torus** At first, we work on the 1-periodic torus  $\mathbb{T} := \mathbb{R}/\mathbb{Z} \cong [0, 1)$  or its multivariate version  $\mathbb{T}^d, d \in \mathbb{N}$ , and with the function spaces of integrable and square integrable, 1-periodic functions,<sup>5</sup>

$$L^1(\mathbb{T}^d) = \left\{ f : \mathbb{T}^d \rightarrow \mathbb{C} \text{ measurable and } \|f\|_{L^1} := \int_{\mathbb{T}^d} |f(x)| dx < \infty \right\} \text{ and}$$

$$L^2(\mathbb{T}^d) = \left\{ f : \mathbb{T}^d \rightarrow \mathbb{C} \text{ measurable and } \|f\|_{L^2} := \left( \int_{\mathbb{T}^d} |f(x)|^2 dx \right)^{1/2} < \infty \right\}.$$

While both are Banach spaces,  $L^2(\mathbb{T}^d)$  is even a Hilbert space with inner product  $\langle f, g \rangle := \int_{\mathbb{T}^d} \overline{f(x)} g(x) dx$ . Moreover, Hölder's inequality yields  $L^2(\mathbb{T}^d) \subset L^1(\mathbb{T}^d)$  on the compact set  $\mathbb{T}^d$ . Additionally, the family of functions  $(e^{2\pi i k x})_{k \in \mathbb{Z}^d}$  forms an orthonormal system in  $L^2(\mathbb{T}^d)$  and this motivates to define the *Fourier coefficients* of any function  $f \in L^1(\mathbb{T}^d)$  by

$$c_k(f) := \langle e^{2\pi i k \cdot}, f \rangle = \int_{\mathbb{T}^d} f(x) e^{-2\pi i k x} dx, \quad k \in \mathbb{Z}^d.$$

If  $f \in L^2(\mathbb{T}^d)$ , the sequence of Fourier coefficients  $c(f)$  satisfies  $c(f) \in \ell^2$  where  $\ell^p := \{v : \mathbb{Z}^d \rightarrow \mathbb{C} : \|v\|_p^p := \sum_{k \in \mathbb{Z}^d} |v_k|^p < \infty\}, p \in (1, \infty)$ , and the function  $f$  admits the representation in terms of the *Fourier series*

$$f(x) = \sum_{k \in \mathbb{Z}^d} c_k(f) e^{2\pi i k x}$$

with convergence of the series in the Hilbert space  $L^2(\mathbb{T}^d)$ . In particular, a function  $f \in L^2(\mathbb{T}^d)$  and its sequence of Fourier coefficients fulfil the *Parseval relation*  $\|f\|_{L^2} = \|c(f)\|_2$ . Beyond that, other types of convergence of the Fourier series like pointwise or uniform convergence are well-studied subjects that we omit at this point.

<sup>5</sup>A good reference for this paragraph is for instance [132, Sec. 1.2 and Sec. 4.1].



**Fourier transform and function spaces on  $\mathbb{R}^d$**  For an integrable function on the real line or more general on  $\mathbb{R}^d$ , i.e. for  $f \in L^1(\mathbb{R}^d)$ , we define the *Fourier transform* as the function  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{C}$  with

$$\hat{f}(v) = \int_{\mathbb{R}^d} f(x) e^{-2\pi i v \cdot x} dx.$$

While the Fourier transform of  $f \in L^1(\mathbb{R}^d)$  is not necessarily integrable,<sup>6</sup> one can extend the Fourier transform to  $L^2(\mathbb{R}^d)$  by a density argument (cf. [132, Sec. 2.2]). On this space of square integrable functions, we have the Parseval relation  $\|f\|_{L^2} = \|\hat{f}\|_{L^2}$  (cf. [132, Thm. 2.22]) and thus the Fourier transform is an isometry of  $L^2(\mathbb{R}^d)$  onto itself. The fact that the inner product is invariant under the Fourier transform, i.e.  $\langle f, g \rangle = \langle \hat{f}, \hat{g} \rangle$  for  $f, g \in L^2(\mathbb{R}^d)$ , is also called *Plancherel's theorem*.

Another space that is preserved under the Fourier transform is the *Schwartz space* denoted by  $\mathcal{S}(\mathbb{R}^d)$ . Its elements are called *Schwartz functions* which are defined as functions  $f : \mathbb{R}^d \rightarrow \mathbb{C}$  such that for all multi-indices  $\alpha, \beta$  the *Schwartz seminorm* satisfies

$$\|f\|_{\alpha, \beta} := \sup_{x \in \mathbb{R}^d} \left| x^\alpha \partial^\beta f(x) \right| \leq C_{\alpha, \beta}$$

for some  $C_{\alpha, \beta} > 0$  (e.g. cf. [61, Def. 2.2.1]).<sup>7</sup> In other words, a function  $f \in \mathcal{S}(\mathbb{R}^d)$  has infinitely many derivatives and not only the function but also its derivatives go to zero faster than any polynomial.

**Lemma 1.2.1.** (Convolution, cf. [132, Thm. 2.1.3]) For  $f, g \in L^1(\mathbb{R}^d)$  their convolution

$$(f * g)(x) = \int_{\mathbb{R}^d} f(y) g(x - y) dy$$

exists for almost every  $x \in \mathbb{R}^d$  and  $f * g \in L^1(\mathbb{R}^d)$  holds.

We summarise the following properties of the Fourier transform.

**Proposition 1.2.2.** (e.g. cf. [132, Sec. 2.1 and Sec. 4.2]) Let  $f, g \in L^1(\mathbb{R}^d)$ . The Fourier transform satisfies

(i) the translation vs. modulation property: For  $x_0, v_0 \in \mathbb{R}^d$

$$\begin{aligned} (f(\cdot - x_0))^\wedge(v) &= e^{-2\pi i x_0 \cdot v} \hat{f}(v), \\ (e^{-2\pi i v_0 \cdot \cdot} f(\cdot))^\wedge(v) &= \hat{f}(v + v_0). \end{aligned}$$

(ii) If for some  $\alpha \in \mathbb{N}^d$  the derivative  $\partial^\alpha f$  exists and  $\partial^\alpha f \in L^1(\mathbb{R}^d)$ , one finds

$$(\partial^\alpha f)^\wedge(v) = (2\pi i v)^\alpha \hat{f}(v).$$

(iii) We have  $(f * g)^\wedge(v) = \hat{f}(v) \cdot \hat{g}(v)$ .

<sup>6</sup>For  $f \in L^1(\mathbb{T}^d)$ , we have  $\hat{f} \in C_0(\mathbb{T}^d)$  which is the space of continuous functions that vanish as  $\|x\|$  goes to infinity.

<sup>7</sup>We use the multi-index notation  $x^\alpha := x_1^{\alpha_1} \cdots x_d^{\alpha_d}$ .

## 1 Preliminaries

Moreover, the Fourier transform admits an inversion formula

$$f(x) = \int_{\mathbb{R}^d} \hat{f}(v) e^{2\pi i v x} dv \text{ for almost every } x \in \mathbb{R}^d$$

if  $f \in L^1(\mathbb{R}^d)$  with  $\hat{f} \in L^1(\mathbb{R}^d)$ .

A connection of the Fourier series for periodic functions and the Fourier transform defined for functions on  $\mathbb{R}^d$  can be made through the *Poisson summation formula*.

**Theorem 1.2.3.** (*Poisson summation formula, e.g. cf. [132, Thm. 4.27]*) Let  $f \in C_0(\mathbb{R}^d)$  be a function satisfying the decay conditions

$$|f(x)| \leq \frac{c}{1 + \|x\|_2^{d+\epsilon}} \quad \text{and} \quad |\hat{f}(v)| \leq \frac{c}{1 + \|v\|_2^{d+\epsilon}}$$

for some  $\epsilon, c > 0$ . Then, the periodisation  $\tilde{f}$  satisfies the Poisson summation formula

$$\tilde{f}(x) = \sum_{l \in \mathbb{Z}^d} f(x + l) = \sum_{k \in \mathbb{Z}^d} \hat{f}(k) e^{2\pi i k x} \quad \text{for all } x \in \mathbb{T}^d,$$

i.e. the  $k$ th Fourier coefficient of  $\tilde{f}$  agrees with the Fourier transform of  $f$  evaluated at  $k$ . Moreover, both series from the Poisson summation formula converge absolutely and uniformly.

While the Poisson summation formula connects Fourier theory on  $\mathbb{T}^d$  and  $\mathbb{R}^d$ , the *sampling theorem* establishes a method to recover a function  $f$  from (infinitely many) function evaluations. But in order to make this possible, the function  $f \in L^2(\mathbb{R}^d)$  needs to be *bandlimited* and this means that its Fourier transform has compact support, i.e.  $\hat{f}(x) = 0$  for almost every  $x$  outside of some compact set. Hence, we define

$$\mathcal{B}_n(\mathbb{R}^d) := \left\{ f \in L^2(\mathbb{R}^d) : \text{supp } \hat{f} \subset [-n, n]^d \right\}$$

as the *space of bandlimited functions* with *bandwidth*  $n$  (cf. [132, p. 86]).

**Theorem 1.2.4.** (*Sampling theorem, Shannon–Whittaker–Kotelnikov, cf. [132, Thm. 2.29]*) Let  $f \in L^1(\mathbb{R}^d) \cap C_0(\mathbb{R}^d) \cap \mathcal{B}_n(\mathbb{R}^d)$  and  $m \geq n > 0$ . Then,  $f$  is uniquely determined by its function evaluations  $f\left(\frac{k}{2m}\right)_{k \in \mathbb{Z}^d}$  with the representation

$$f(x) = \sum_{k \in \mathbb{Z}^d} f\left(\frac{k}{2m}\right) \text{sinc}(2mx - k),$$

converging in  $L^2(\mathbb{R}^d)$  as well as absolutely and uniformly on  $\mathbb{R}^d$ . Here, the function  $\text{sinc}$  denotes the *sinus cardinalis*,

$$\text{sinc}(x) = \frac{\sin(\pi x)}{\pi x} \quad \text{for } x \in \mathbb{R},$$

or the tensor products with itself in higher dimensions.

*Proof.* The univariate case is classic and for instance it is presented in [132]. Based on this, the multivariate extension with the presented tensor product structure can be deduced directly (cf. [61, Thm. 6.6.9]).  $\square$

The following lemma is a classical exercise in a course on Fourier analysis.

**Lemma 1.2.5.** *The functions  $s_k \in \mathcal{B}_n(\mathbb{R})$ ,  $s_k(x) := \text{sinc}(2nx - k)$ ,  $k \in \mathbb{Z}$ , are orthogonal. For  $f \in L^1(\mathbb{R}^d) \cap C_0(\mathbb{R}^d) \cap \mathcal{B}_n(\mathbb{R}^d)$ ,  $m \geq n$  we have*

$$\|f\|_2^2 = \frac{1}{(2m)^d} \sum_{k \in \mathbb{Z}^d} \left| f\left(\frac{k}{2m}\right) \right|^2.$$

*Proof.* The orthogonality of the translated sinc functions  $s_k$  is a direct consequence of Plancherel's theorem since  $(\text{sinc})^\wedge = \mathbb{1}_{[-1/2, 1/2]}$  is the *indicator function* on  $[-1/2, 1/2]$  and

$$\begin{aligned} \langle s_k, s_j \rangle &= \langle \hat{s}_k, \hat{s}_j \rangle = \int_{\mathbb{R}^2} e^{-2\pi i k v / (2n)} \frac{1}{(2n)^2} \mathbb{1}_{[-n, n]}(v) e^{2\pi i j v / (2n)} \mathbb{1}_{[n, n]}(v) dv \\ &= \int_{-n}^n \frac{1}{(2n)^2} e^{2\pi i (j-k)v / (2n)} dv \\ &= \frac{1}{2n} \delta_{k,j} \end{aligned}$$

for  $k, j \in \mathbb{Z}$ . Analogously, a similar result can be given for the multivariate version of sinc functions. For the second statement, we can use the sampling theorem in  $L^2(\mathbb{R}^d)$  and conclude

$$\left\| \|f\|_2 - \left\| \sum_{\|k\|_\infty \leq l} f\left(\frac{k}{2m}\right) s_k \right\|_2 \right\| \leq \left\| f - \sum_{\|k\|_\infty \leq l} f\left(\frac{k}{2m}\right) s_k \right\|_2 \xrightarrow{l \rightarrow \infty} 0.$$

This can be used to establish the statement

$$\begin{aligned} \|f\|_2^2 &= \lim_{l \rightarrow \infty} \left\| \sum_{\|k\|_\infty \leq l} f\left(\frac{k}{2m}\right) s_k \right\|_2^2 \\ &= \lim_{l \rightarrow \infty} \sum_{\|k\|_\infty \leq l} \left\| f\left(\frac{k}{2m}\right) s_k \right\|_2^2 \\ &= \sum_{k \in \mathbb{Z}^d} \left| f\left(\frac{k}{2m}\right) \right|^2 \|s_k\|_2^2 = \frac{1}{(2m)^d} \sum_{k \in \mathbb{Z}^d} \left| f\left(\frac{k}{2m}\right) \right|^2, \end{aligned}$$

where we used the orthogonality of  $s_k, k \in \mathbb{Z}^d$ , in the second equality.  $\square$

**Discrete Fourier transform** An obvious way to approximate the Fourier coefficients of a 1-periodic function from a finite number of evaluations is to compute the *discrete Fourier transform* (DFT)  $\hat{f} \in \mathbb{C}^K$ , where

$$\hat{f}_j = \sum_{k=0}^{K-1} f_k e^{-2\pi i k j / K}, \quad j = 0, \dots, K-1$$

for some vector  $f \in \mathbb{C}^K$ ,  $K \in \mathbb{N}$ . One of the main advantages of the DFT is that it can be implemented in a fast manner by the *fast Fourier transform* (FFT) just using  $\mathcal{O}(K \log K)$  operations. The following lemma shows how accurate this approximation is.

**Lemma 1.2.6.** (Alias lemma, e.g. cf. [132, Thm. 3.3]) Suppose that  $f \in C(\mathbb{T})$ , i.e. the function  $f$  is 1-periodic and continuous. Moreover, let the sequence of Fourier coefficients  $c(f)$  be absolutely summable. If we then take a vector  $\tilde{f} \in \mathbb{C}^K$ ,  $\tilde{f}_k = f\left(\frac{k}{K}\right)$ ,  $k = 0, \dots, K-1$ , the discrete Fourier transform of  $\tilde{f}$  satisfies the alias formula

$$\frac{1}{K} \hat{\tilde{f}}_k = \sum_{\ell \in \mathbb{Z}} c_{k+\ell K}(f), \quad k \in \mathbb{Z}.$$

We remark that this can be generalised to higher dimensions. Furthermore, the periodisation of a bandlimited function has Fourier coefficients given by the Fourier transform evaluated at integers due to the Poisson summation formula. By the limited support of the Fourier transform and the alias formula, the DFT then computes the first Fourier coefficients exactly in this case.

### 1.3 Some auxiliary functions

In this section, we introduce some auxiliary functions which are used throughout the thesis.

**Bessel functions** While Bessel functions are usually motivated through the analysis of certain differential equations, they naturally appear in our work when we study *radial functions* and their Fourier transform.<sup>8</sup>

**Definition 1.3.1.** (Radial function, e.g. cf. [132, Sec. 4.2.4]) A function  $f : \mathbb{R}^d \rightarrow \mathbb{C}$  is *radial* if there exists a function  $\tilde{f} : [0, \infty) \rightarrow \mathbb{C}$  such that  $f(x) = \tilde{f}(\|x\|_2)$ .

It is straightforward to show that a radial  $L^1$ -function has also a radial Fourier transform (cf. [132, Cor. 4.30]). Furthermore, the Fourier transform of a  $d$ -variate radial function  $f$ ,  $d > 1$ , can be expressed as an one-dimensional integral

$$\hat{f}(v) = 2\pi \|v\|_2^{-d/2+1} \int_0^\infty \tilde{f}(r) J_{d/2-1}(2\pi r \|v\|_2) r^{d/2} dr, \quad (1.1)$$

where for  $x \geq 0$

$$J_\nu(x) = \sum_{k=0}^\infty \frac{(-1)^k}{k! \Gamma(k + \nu + 1)} \left(\frac{x}{2}\right)^{2k+\nu} \quad (1.2)$$

denotes the *Bessel function (of the first kind) of order  $\nu > -1$* , see [61, p. 577]. Therein,  $\Gamma$  denotes the Gamma function. Due to its importance for computing radial Fourier transforms, the integral on the right hand side of (1.1) is also called *Hankel transform* of order  $d/2 - 1$  of  $\tilde{f}$ . The graphs of  $J_\nu$  for a few values of  $\nu$  presented in Figure 1.1 show an oscillating behaviour of these Bessel functions on the positive real line. We include various other properties of Bessel functions in the following lemma.

**Lemma 1.3.2.** Let  $j_{\nu,k}$  be the  $k$ th smallest positive zero of  $J_\nu$  for  $\nu > -1$ . Then, we have

- (i)  $0 < j_{\nu,1} < j_{\nu+1,1} < j_{\nu,2}$ ,
- (ii)  $j_{\nu,1} = \nu + 1.855757\nu^{1/3} + \mathcal{O}(\nu^{-1/3})$  as  $\nu \rightarrow \infty$  and

---

<sup>8</sup>A good reference for an overview over Bessel functions is [154].

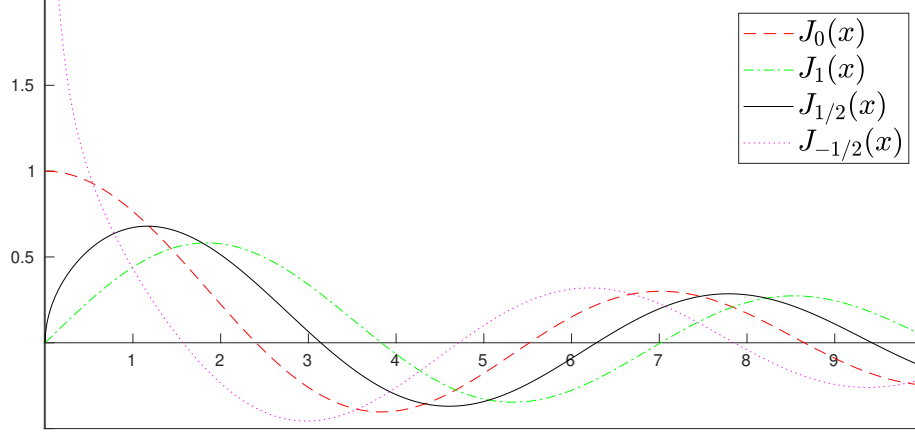


Figure 1.1: The graphs of some Bessel functions. We included a Bessel function with negative order which is possible by extending (1.2) to  $\nu \in \mathbb{C}$  and  $x \in \mathbb{C}$  where  $(\frac{x}{2})^\nu$  and  $\Gamma(\nu + 1)$  can be defined. The latter holds if  $\nu > -1$  and the real part of  $x$  satisfies  $\Re(x) > 0$ . Here,  $J_{-1/2}$  has a pole at  $x = 0$ .

(iii)  $j_{-1/2,1} = \frac{\pi}{2}$ ,  $j_{0,1} \approx 2.4048$ ,  $j_{1/2,1} = \pi$  and  $j_{1,1} \approx 3.8317$ .

Moreover, the Bessel functions satisfy the recurrence relation

$$J_{\nu+1}(x) = \frac{2\nu}{x} J_\nu(x) - J_{\nu-1}(x), \quad x > 0$$

and the asymptotic expansion

$$J_\nu(x) = \left( \left( \frac{2}{\pi x} \right)^{1/2} \cos \left( x - \frac{1}{2} \nu \pi - \frac{1}{4} \pi \right) \right) + \mathcal{O}(x^{-3/2}) \quad \text{as } x \rightarrow \infty.$$

Finally, Bessel functions can be expressed by trigonometric functions in the special cases

$$J_{-1/2}(x) = \left( \frac{2}{\pi x} \right)^{1/2} \cos(x) \quad \text{and} \quad J_{1/2}(x) = \left( \frac{2}{\pi x} \right)^{1/2} \sin(x).$$

*Proof.* For properties (i)-(iii) see [154, Chap. XV]. The representations of  $J_{\pm 1/2}$  can be found in [154, p. 54] while the recurrence formula is for instance given in [154, p. 45]. A reference for the asymptotic formula is [154, p. 199].  $\square$

**Extremal minorants and localising functions** Due to the *uncertainty principle*, it is not possible to find functions which are compactly supported in spatial and frequency domain. Here, we introduce functions that circumvent this issue by being localised in one domain through a compact support while they are also localised in the other domain as a minorant or majorant of a certain function. In [13] as cited in [150] or in [14] as quoted in [7, p. 116] respectively, Beurling introduced the function<sup>9</sup>

$$B : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto B(x) = \frac{\sin^2(\pi x)}{\pi^2} \left[ \sum_{k \in \mathbb{Z}} \frac{\operatorname{sgn}(k)}{(x-k)^2} + \frac{2}{x} \right], \quad (1.3)$$

<sup>9</sup>For  $l \in \mathbb{Z}$ , the value of  $B$  should be understood in the sense of a limit, i.e.  $B(l) := \lim_{x \rightarrow l} B(x) = \operatorname{sgn}(l)$ .

## 1 Preliminaries

where the *sign function* is defined as

$$\operatorname{sgn} : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto \operatorname{sgn}(x) = \begin{cases} 1, & x \geq 0 \\ -1, & x < 0. \end{cases}$$

While  $B$  is graphically displayed in Figure 1.2 (a), we summarise some of its properties.

**Proposition 1.3.3.** ([150, Thm. 8]) *The Beurling function  $B$*

- (i) *majorises the sign function, i.e.  $\operatorname{sgn}(x) \leq B(x)$ ,*
  - (ii) *can be extended to an entire function of exponential type<sup>10</sup>  $2\pi$ ,*
  - (iii) *has a distributional<sup>11</sup> Fourier transform  $\hat{B}$  with  $\operatorname{supp} \hat{B} = [-1, 1]$  and*
  - (iv) *its difference to the sign function measured in  $L^1$  satisfies  $\int_{\mathbb{R}} B(x) - \operatorname{sgn}(x) dx = 1$ .*
- Furthermore, all functions  $f$  satisfying (i)-(iii) automatically fulfil

$$\int_{\mathbb{R}} f(x) - \operatorname{sgn}(x) dx \geq 1 \tag{1.4}$$

with equality if and only if  $f = B$ .

The last part of Proposition 1.3.3 motivates to speak about  $B$  as an *extremal majorant* of the sign function. On the other hand, one directly obtains an *extremal minorant* of  $\operatorname{sgn}$  by using  $-B(-x)$  (cf. [150]). As it is of particular interest to find minorant and majorants for the indicator function  $\mathbb{1}_{[-1,1]}$ , e.g. for applications in control theory [83] or in order to bound the smallest and largest singular values of Vandermonde matrices (cf. [117, 6]), Selberg [146] as cited in [7, 6] introduced the functions

$$\begin{aligned} s_+(x) &= \frac{1}{2} (B(x+1) + B(1-x)) \quad \text{and} \\ s_-(x) &= -\frac{1}{2} (B(-1-x) + B(x-1)) = \frac{\sin^2(\pi x)}{\pi^2} \left( \frac{2}{1-x^2} + \frac{1}{x^2} \right). \end{aligned}$$

As a consequence of Proposition 1.3.3 we know that  $s_-, s_+$  fulfil variants of (i)-(iv) in Proposition 1.3.3, i.e. they have exponential type  $2\pi$  and admit  $s_- \leq \mathbb{1}_{[-1,1]} \leq s_+$ ,  $\operatorname{supp} \hat{s}_- = \operatorname{supp} \hat{s}_+ = [-1, 1]$  as well as

$$\int_{\mathbb{R}} \mathbb{1}_{[-1,1]}(x) - s_-(x) dx = \int_{\mathbb{R}} s_+(x) - \mathbb{1}_{[-1,1]}(x) dx = 1. \tag{1.5}$$

Moreover, any function  $f$  with exponential type  $2\pi$  being either a majorant with  $\mathbb{1}_{[-1,1]} \leq f$  or a minorant with  $\mathbb{1}_{[-1,1]} \geq f$  has the property that

$$\int_{\mathbb{R}} |f(x) - \mathbb{1}_{[-1,1]}(x)| dx \geq 1,$$

<sup>10</sup>A function  $f : \mathbb{C} \rightarrow \mathbb{C}$  is called *entire* if it is holomorphic everywhere. Additionally, having an entire function  $f$  of *exponential type*  $\sigma > 0$  means that for every  $\epsilon > 0$  there exists  $c_\epsilon > 0$  such that  $|f(z)| \leq c_\epsilon e^{(\sigma+\epsilon)|z|}$  for all  $z \in \mathbb{C}$  (cf. [150, p. 186]). Properties (ii) and (iii) in Proposition 1.3.3 are connected through the Paley-Wiener-Schwartz theorem (e.g. cf. [138, Thm. 19.3] or [72, Sec. 7.3]).

<sup>11</sup>Even though  $B \notin L^p(\mathbb{R})$  for all  $p \in [1, \infty)$ , one can define its Fourier transform as a *distribution* via  $\langle \hat{B}, \phi \rangle := \langle B, \hat{\phi} \rangle$  for all test functions  $\phi \in \mathcal{S}(\mathbb{R})$ .

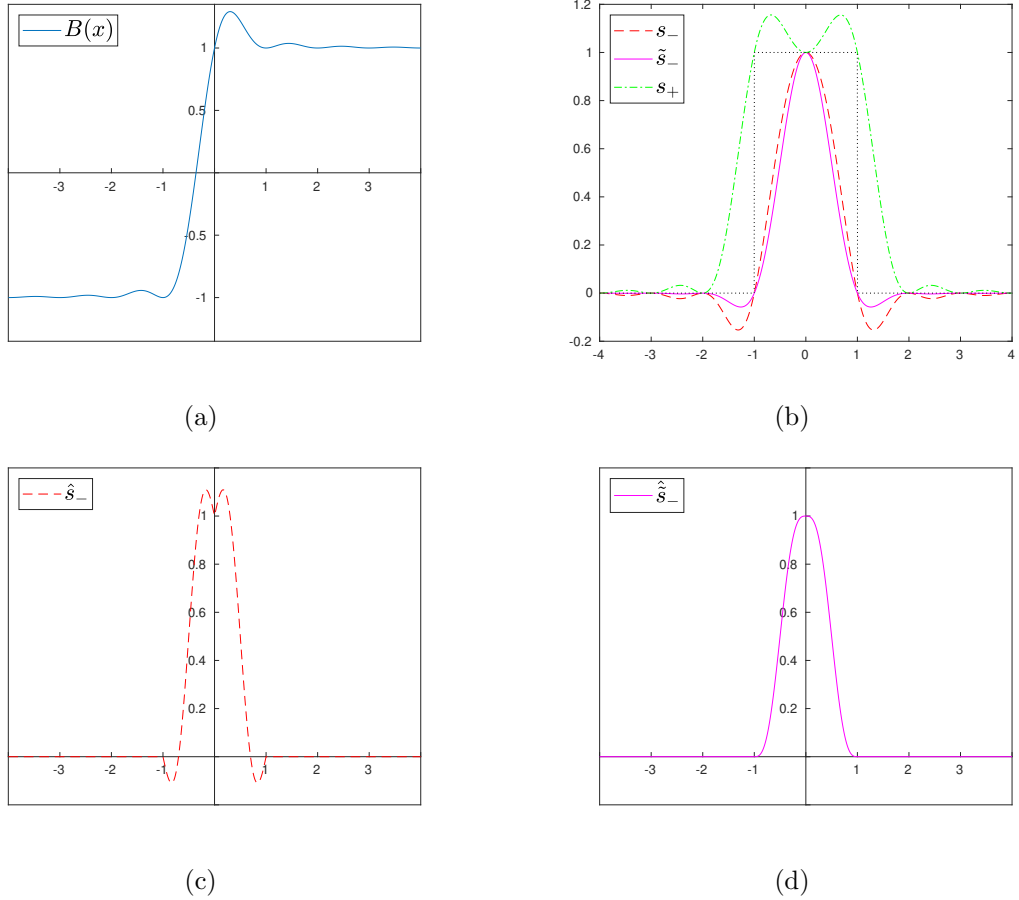


Figure 1.2: Extremal Functions of Beurling and Selberg compared to construction of Diederichs. The Beurling function  $B$  is an extremal majorant to the sign function (a) where we used the first 200 summands in (1.3) for the visualisation. Due to Selberg, it can be used to establish a minorant and a majorant to the characteristic function of the interval  $[-1, 1]$  (b). But the Fourier transform of Selberg's minorant is not maximal in zero (c) while a construction by Diederichs fulfills this, see (b) and (d).

see [150, p. 186] or [37].<sup>12</sup> Hence,  $s_-, s_+$  are an extremal minorant and majorant of  $\mathbb{1}_{[-1,1]}$ . However, they are not unique in contrast to the situation for the sign function where the Beurling function is the only extremal majorant, see (1.4). This can be seen by a construction of Diederichs [37, Prop. 2.13] who worked with

$$\tilde{s}_-(x) = \frac{\sin^2(\pi x)}{\pi^2} \left( \frac{1}{1-x^2} + \frac{1}{x^2} \right)$$

which is another extremal minorant of  $\mathbb{1}_{[-1,1]}$ .<sup>13</sup> Both minorants and the majorant of Selberg are depicted in Figure 1.2 (b). Compared to the minorant  $s_-$  of Selberg, the

<sup>12</sup>We remark that the case of a indicator function on a interval  $[a, b]$  with  $b - a \notin \mathbb{Z}$  behaves different, e.g. see [102].

<sup>13</sup>One can see that easily by computing the Fourier transforms of  $s_-, \tilde{s}_-$  and observing that their values in zero agree which means that  $\int_{\mathbb{R}} s_-(x) dx = \int_{\mathbb{R}} \tilde{s}_-(x) dx$ .

## 1 Preliminaries

minorant  $\tilde{s}_-$  has the advantage that its Fourier transform  $\hat{\tilde{s}}_-$  is maximal in zero whereas  $\hat{s}_-$  has even a local minimum in zero. To see this, one can use [37, Lem. 2.21] in order to obtain

$$\hat{s}_-(v) = \left( \frac{\sin(2\pi|v|)}{\pi} + 1 - |v| \right) \mathbb{1}_{[-1,1]}(v) \quad \text{and} \quad \hat{\tilde{s}}_-(v) = \left( \frac{\sin(2\pi|v|)}{2\pi} + 1 - |v| \right) \mathbb{1}_{[-1,1]}(v).$$

Hence, it is trivial to analyse the extremum for  $v = 0$ . Additionally, the different types of extrema can be seen in their graphs, see Figure 1.2 (c) and (d). We remark that in [37, 38] Diederichs also coined the term *localising function* for minorants of the indicator function of an interval with compactly supported Fourier transform.

For higher dimensions, one has to concretise the multivariate generalisation of  $\mathbb{1}_{[-1,1]}$  at first. An obvious choice is to consider the box function  $\mathbb{1}_{[-1,1]^d}$  and this was done in [24, 23] where the authors construct minorants for small dimensions  $d \leq 5$  while they remark that “in higher dimensions no extremal results for the box minorant problem are known” (cf. [23, p. 2]). Compared to the box-minorant problem, the question whether there are functions  $f_-, f_+$  such that

$$f_-(x) \leq \mathbb{1}_{B_n(0)} \leq f_+(x) \quad \text{and} \quad \text{supp } \hat{f}_-, \text{supp } \hat{f}_+ \subset B_q(0), \quad (1.6)$$

where  $B_q(0) = \{x \in \mathbb{R}^d : \|x\|_2 \leq q\}$  is the Euclidean ball, was addressed in [71]. Since the radial structure of the problem allows to transform it back to a univariate problem, the problem (1.6) is simpler than the box problem and thus extremal minorants and majorants can be found for any  $n, q$ , see [71]. However, the work by Goncalves [59] shows that a localising function  $f_-$  with  $\hat{f}_-(0) > 0$  can only exist if  $n \cdot q$  is larger than a certain dimension dependent critical radius. Both works [71] and [59] have the drawback that the localising functions are not given explicitly. An explicit construction for a function satisfying (1.6) is given in [83] by taking

$$f_-(x) = 4\pi^2(n^2 - x^2)|\hat{\varphi}(x)|^2 \quad \text{or} \quad \hat{f}_-(v) = (4\pi^2n^2 + \Delta)(\varphi * \varphi)(v), \quad (1.7)$$

where  $\Delta = \sum_{s=1}^d \frac{\partial^2}{\partial x_s^2}$  is the Laplace operator and  $\varphi$  a function with  $\text{supp } \varphi \subset B_{q/2}(0)$ . This idea was used in a modified form for instance in [86] to prove inequalities for singular values of Vandermonde matrices. In Subsection 2.2.1 of this work, we tune the function  $\varphi$  such that  $\hat{f}_-$  becomes maximal in zero as this will be an important ingredient in order to derive Subsection 2.2.2. This need emerged already in [53] studying the problem in the univariate case and in [37, 38] where not only the univariate but also the bivariate case is addressed. Our approaches extend to arbitrary dimensions and allow to make  $n \cdot q$  as small as possible.

**Trigonometric polynomials and kernels** In Chapter 3, we approximate and interpolate measures by *trigonometric polynomials*, i.e. functions  $p : \mathbb{T}^d \rightarrow \mathbb{C}$ ,

$$x \mapsto p(x) = \sum_{k \in \mathcal{K}} \mathbf{p}_k e^{2\pi i k x}$$

for some finite set  $\mathcal{K} \subset \mathbb{Z}^d$  and *coefficients*  $\mathbf{p}_k = c_k(p) \in \mathbb{C}$ ,  $k \in \mathbb{Z}^d$ . In one dimension, we typically have the symmetric case with  $\mathcal{K} = \{-n, \dots, n\}$  for some  $n \in \mathbb{N}$  called *degree* of  $p$ . The latter is denoted by  $\text{deg}(p) = n$ . As already for minorants in the previous paragraph, simple extensions to higher dimensions include the box case, i.e.  $\mathcal{K} = \{-n, \dots, n\}^d$ , or the



radial case with  $\mathcal{K} = \{k \in \mathbb{Z}^d : \|k\|_2^2 \leq n^2\}$ . In order to distinguish the two, we denote the *maximal degree*  $n$  in the box case by  $\deg_\infty(p) = n$  whereas the *Euclidean degree*  $\deg_2(p) = n$  is used in the radial case. Moreover, we define the space of *polynomials with maximal degree*  $n$  as

$$\mathcal{P}^{n,d,\infty} = \left\{ p : \mathbb{T}^d \rightarrow \mathbb{C}, x \mapsto p(x) = \sum_{k \in \{-n, \dots, n\}^d} \mathbf{p}_k e^{2\pi i k x} \right\}$$

and the space of *radial polynomials with Euclidean degree*  $n$  as

$$\mathcal{P}^{n,d,2} = \left\{ p : \mathbb{T}^d \rightarrow \mathbb{C}, x \mapsto p(x) = \sum_{k: \|k\|_2^2 \leq n^2} \mathbf{p}_k e^{2\pi i k x} \right\}.$$

In the course of this work, we utilise some well-established trigonometric polynomials, e.g. cf. [132]. Following the notation in [120], we introduce the *Dirichlet kernel* and its *modified* version.

**Definition 1.3.4** ((Modified) Dirichlet kernel). The *Dirichlet kernel*  $D_n \in \mathcal{P}^{n,d,\infty}$  is the trigonometric polynomial with

$$D_n(x) := \sum_{k=-n}^n e^{2\pi i k x} = \begin{cases} \frac{\sin((2n+1)\pi x)}{\sin(\pi x)}, & x \neq 0, \\ 2n+1, & x = 0, \end{cases} \quad (1.8)$$

in the univariate case and its tensor product  $D_n(x_1, \dots, x_d) := D_n(x_1) \cdots D_n(x_d)$  for the multivariate situation. The second equality in (1.8) follows directly by the geometric sum formula. In contrast to that, the *modified Dirichlet kernel*

$$d_n(x) := \sum_{k=0}^n e^{2\pi i k x} = \begin{cases} e^{\pi i n x} \frac{\sin((n+1)\pi x)}{\sin(\pi x)}, & x \neq 0, \\ n+1, & x = 0, \end{cases}$$

only includes the nonnegative frequencies. Again, we follow a tensor product approach by setting  $d_n(x_1, \dots, x_d) := d_n(x_1) \cdots d_n(x_d)$  in higher dimensions.

The Dirichlet kernel appears naturally since every truncated Fourier series can be understood as a convolution of the original function  $f \in L^2(\mathbb{T}^d)$  with the Dirichlet kernel. Hence, convolutions with the Dirichlet kernel are well-studied in harmonic analysis in order to describe convergence of Fourier series. Taking the squared absolute value of  $d_n$  still yields a polynomial from  $\mathcal{P}^{n,d,\infty}$ . Properly normalised, this polynomial is known as the *Fejér kernel*, see [132, Exa. 1.15].

**Definition 1.3.5** (Fejér kernel). The *Fejér kernel*  $F_n \in \mathcal{P}^{n,d,\infty}$  is the trigonometric polynomial defined as

$$F_n(x) := \frac{1}{(n+1)^d} |d_n(x)|^2 = \sum_{k \in \{-n, \dots, n\}^d} \left( \prod_{s=1}^d 1 - \frac{k_s}{n+1} \right) e^{2\pi i k x}.$$

In contrast to the Dirichlet kernel, the Fejér kernel is a *summation kernel* yielding uniform convergence of the convolution of the Fejér kernel with any continuous function  $f \in C(\mathbb{T}^d)$  to this function  $f$ , see [132, Thm. 1.17]. However, convolution with the Fejér kernel does not lead to the optimal convergence rate among all approximations through a convolution. This is then achieved by the *Jackson kernel*, see [77, pp. 2 ff.].

**Definition 1.3.6** (Jackson kernel). The *Jackson kernel*

$$J_{2m-2}(x) = \frac{3}{m(2m^2+1)} \frac{\sin^4(m\pi x)}{\sin^4(\pi x)}, \quad m \in \mathbb{N},$$

has degree  $n = 2m - 2$  and is easily extended to higher dimensions by the product approach with  $J_n(x_1, \dots, x_d) := J_n(x_1) \cdots J_n(x_d) \in \mathcal{P}^{n,d,\infty}$ .

Finally, a natural element of  $\mathcal{P}^{n,d,2}$  is the *radial Dirichlet kernel*

$$D_{\text{rad},n}(x) := \sum_{k \in \mathbb{Z}^d: \|k\|_2 \leq n} e^{2\pi i k x}.$$

For this kernel, we can derive a tail estimate analogously to the classical analysis of  $D_n$ . We will use this bound in Chapter 3.

**Lemma 1.3.7.** For  $x \in \mathbb{T}^d \setminus \{0\}$  there is a constant  $c_d > 0$  such that

$$|D_{\text{rad},n}(x)| \leq \frac{c_d n^{d-1}}{\min_{j \in \mathbb{Z}^d} |x + j|_\infty}.$$

*Proof.* Assume without loss of generality  $|x_1|_{\mathbb{T}} := \min_{j_1 \in \mathbb{Z}} |x_1 + j_1| = \min_{j \in \mathbb{Z}^d} |x + j|_\infty$  and set  $x = (x_1, x')$  for  $x' \in \mathbb{T}^{d-1}$ . Denoting  $e_1 = (1, 0, \dots, 0)^\top \in \mathbb{N}^d$ , we have

$$\begin{aligned} & (1 - e^{2\pi i x_1}) D_{\text{rad},n}(x) \\ &= \sum_{k \in \mathbb{Z}^d: \|k\|_2 \leq n} e^{2\pi i k x} - \sum_{k \in \mathbb{Z}^d: \|k\|_2 \leq n} e^{2\pi i (k + e_1) x} \\ &= \sum_{k' \in \mathbb{Z}^{d-1}: \|k'\|_2 \leq n} e^{2\pi i k' x'} \left( \sum_{\substack{k_1 \geq -\sqrt{n^2 - \|k'\|_2^2} \\ k_1 < -\sqrt{n^2 - \|k'\|_2^2} + 1}} e^{2\pi i k_1 x_1} - \sum_{\substack{k_1 > \sqrt{n^2 - \|k'\|_2^2} \\ k_1 \leq \sqrt{n^2 - \|k'\|_2^2} + 1}} e^{2\pi i k_1 x_1} \right) \\ &= \sum_{k' \in \mathbb{Z}^{d-1}: \|k'\|_2 \leq n} e^{2\pi i k' x'} \left( e^{-2\pi i \lfloor \sqrt{n^2 - \|k'\|_2^2} \rfloor x_1} - e^{2\pi i \lfloor \sqrt{n^2 - \|k'\|_2^2} \rfloor + 1 x_1} \right) \end{aligned}$$

where we use the rounding operation  $\lfloor y \rfloor = k \in \mathbb{Z}$  if  $y \in [k, k + 1)$ . Taking the absolute value on both sides gives

$$|D_{\text{rad},n}(x)| \leq \frac{2 \sum_{k' \in \mathbb{Z}^{d-1}: \|k'\|_2 \leq n} 1}{|\sin(\pi x_1)|} \leq \frac{c_d n^{d-1}}{\min_{j \in \mathbb{Z}^d} |x + j|_\infty}$$

for some constant  $c_d > 0$  which can be bounded by  $2^{d-1}$  as  $\{k' \in \mathbb{Z}^{d-1} : \|k'\|_2 \leq n\}$  is a subset of  $[-n, n]^{d-1}$ .  $\square$

## 1.4 Wasserstein metric and optimal transport

The field of *optimal transport* allows to define the *Wasserstein distance* as a metric for measures.<sup>14</sup> Classically, the Wasserstein distance compares a measure  $\mu$  on some space  $\mathcal{X}$

<sup>14</sup>For example, see [129, 151, 141] for an overview on optimal transport.

to  $\nu$  on  $\mathcal{Y}$  where  $\mathcal{X}, \mathcal{Y}$  are typically complete and separable metric spaces equipped with their Borel  $\sigma$ -algebra.<sup>15</sup> For any of these Borel measures  $\mu$  on  $\mathcal{X}$ , we can define the *total variation*  $|\mu| \in \mathcal{M}_+(\mathcal{X})$  by

$$|\mu|(A) = \sup \left\{ \sum_i |\mu(A_i)| : A = \bigcup_i A_i, A_i \cap A_j = \emptyset \text{ if } i \neq j \right\}$$

and say that  $\mu \in \mathcal{M}(\mathcal{X})$  has finite total variation if  $\|\mu\|_{\text{TV}} := |\mu|(\mathcal{X}) < \infty$  (cf. [141, p. 117]). The *space of complex Borel measures with finite total variation* on  $\mathcal{X}$  is denoted by  $\mathcal{M}(\mathcal{X})$  but in the first part of this section we focus on *probability measures*  $\mu \in \mathcal{M}_{+,1}(\mathcal{X})$ ,  $\nu \in \mathcal{M}_{+,1}(\mathcal{Y})$ , i.e. nonnegative measures with  $\mu(\mathcal{X}) = \nu(\mathcal{Y}) = 1$ .<sup>16</sup> For a sequence  $(\mu_k)_{k \in \mathbb{N}} \in \mathcal{M}(\mathcal{X})$ , we say that it *converges weakly* to some  $\mu \in \mathcal{M}(\mathcal{X})$  if

$$\int_{\mathcal{X}} f(x) d\mu_k(x) \rightarrow \int_{\mathcal{X}} f(x) d\mu(x) \quad (1.9)$$

as  $k \rightarrow \infty$  for all test functions  $f \in C(\mathcal{X}) \cap L^\infty(\mathcal{X})$ , i.e. all continuous and bounded functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  (cf. [151, p. 96]). We then write  $\mu_k \rightharpoonup \mu$ .

### Monge vs. Kantorovich formulation and Wasserstein distance for probability measures

Transforming one measure  $\mu \in \mathcal{M}_{+,1}(\mathcal{X})$  into another  $\nu \in \mathcal{M}_{+,1}(\mathcal{Y})$  can be modelled by the *push-forward measure*  $T_{\#}\mu \in \mathcal{M}_{+,1}(\mathcal{Y})$  where  $T : \mathcal{X} \rightarrow \mathcal{Y}$  is a measurable map and

$$T_{\#}\mu(A) = \mu(T^{-1}(A)) = \mu(\{x \in \mathcal{X} : T(x) \in A\})$$

for any Borel set  $A \subset \mathcal{Y}$ . This allows to formulate the *Monge problem*.

**Definition 1.4.1.** (Monge formulation, cf. [129, Rem. 2.7]) For measures  $\mu \in \mathcal{M}_{+,1}(\mathcal{X})$ ,  $\nu \in \mathcal{M}_{+,1}(\mathcal{Y})$  and a cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  the *Monge cost* between  $\mu$  and  $\nu$  is

$$\min_T \int_{\mathcal{X}} c(x, T(x)) d\mu(x) \quad \text{s.t.} \quad \nu = T_{\#}\mu.$$

Every admissible push forward map  $T$  is called *transport map*.

As there are situations where the Monge problem is not solvable, e.g. if one considers two discrete measures with not compatible number of support points (cf. [129, pp. 10-11]), one typically considers the relaxation to the *Kantorovich problem*.

**Definition 1.4.2.** (Kantorovich, e.g. [129, Rem. 2.13]) Let  $\Pi(\mu, \nu)$  denote the *space of all coupling measures* of  $\mu$  and  $\nu$  given by all measures  $\pi \in \mathcal{M}_{+,1}(\mathcal{X} \times \mathcal{Y})$  such that their *marginals* are  $\mu$  and  $\nu$  respectively. The Kantorovich formulation of optimal transport is then

$$\min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y).$$

<sup>15</sup>Spaces like  $\mathcal{X}, \mathcal{Y}$  are called *Polish spaces*, see [151, p. XX]. We usually consider the torus  $\mathbb{T}^d$  or possibly compact subsets of  $\mathbb{R}^d$  for  $\mathcal{X}, \mathcal{Y}$ .

<sup>16</sup>We denote the space of nonnegative measures on some Polish space  $\mathcal{X}$  by  $\mathcal{M}_+(\mathcal{X})$ . Additionally, the space of real valued but signed measures is denoted by  $\mathcal{M}_{\mathbb{R}}(\mathcal{X})$ .

In the literature, there exists a well-established theory about existence and uniqueness of minimisers. Furthermore, there are results that present sufficient conditions for the equivalence of both approaches, see [129, 141, 151] and the references therein. In special cases, e.g. when both measures are discrete or univariate or one is discrete and the other has a density with respect to the Lebesgue measure (*semidiscrete optimal transport*), a richer theory is known, e.g. see [129, Sec. 2.6 and Chapter 5]. For probability measures on a metric space  $\mathcal{X} = \mathcal{Y}$  one can then define the *Wasserstein distance*.

**Definition 1.4.3.** (Wasserstein distance on  $\mathcal{M}_{+,1}(\mathcal{X})$ , e.g. cf. [129, p. 21]) Let  $(\mathcal{X}, d)$  be a metric space and  $p \in [1, \infty)$ . Then, the  $p$ -Wasserstein distance of  $\mu, \nu \in \mathcal{M}_{+,1}(\mathcal{X})$  is given by

$$W_p(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \left( \int_{\mathcal{X} \times \mathcal{X}} d(x, y)^p d\pi(x, y) \right)^{1/p}$$

as the Kantorovich cost between  $\mu$  and  $\nu$  with cost function  $d(\cdot, \cdot)^p$  defines a metric on  $\mathcal{M}_{+,1}(\mathcal{X})$ .

Finally, we remark that weak convergence  $\mu_k \rightharpoonup \mu$  is then equivalent to  $W_p(\mu_k, \mu) \rightarrow 0$  as  $k \rightarrow \infty$ , see [151, Thm. 6.9]. This enables to quantify a convergence rate for weakly convergent sequences.

**Dual formulation and extension of Wasserstein to complex measures** Extensions of the Wasserstein distance to arbitrary complex measures can be made by the dual formulation of the Kantorovich problem.

**Theorem 1.4.4.** (*Kantorovich duality*, [151, Thm. 5.10]) Let  $\mu, \nu \in \mathcal{M}_{+,1}(\mathcal{X})$ . If the cost function is lower semi-continuous, one has the duality

$$\min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) = \sup_{\substack{f \in C(\mathcal{X}) \cap L^\infty(\mathcal{X}), g \in C(\mathcal{Y}) \cap L^\infty(\mathcal{Y}): \\ f(x) + g(y) \leq c(x, y)}} \int_{\mathcal{X}} f(x) d\mu(x) + \int_{\mathcal{Y}} g(y) d\nu(y)$$

where  $L^\infty(\mathcal{X})$  denotes the set of bounded functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  and  $C(\mathcal{X})$  contains all those functions  $f$  which are continuous. In the case  $\mathcal{X} = \mathcal{Y}$  and  $c(x, y) = d(x, y)$  for some metric  $d$  on  $\mathcal{X}$ , the Kantorovich-Rubenstein formula, see [151, p. 60],

$$\min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} d(x, y) d\pi(x, y) = \sup_{f: \text{Lip}(f) \leq 1} \int_{\mathcal{X}} f(x) (d\mu(x) - d\nu(x))$$

holds where the supremum on the right hand side is over all Lipschitz continuous functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  with  $\text{Lip}(f) := \sup_{x \neq y} \frac{|f(x) - f(y)|}{d(x, y)} \leq 1$ .

Hence, the 1-Wasserstein distance has a very simple dual formulation in terms of Lipschitz functions and we now use this to generalize the concept of Wasserstein distances to complex valued measures from  $\mathcal{M}(\mathcal{X})$  which additionally need not to have equal mass if  $\mathcal{X}$  is bounded.<sup>17</sup> Even if [129, p. 98] mentions that the Wasserstein distance can be extended to signed measures with equal mass through the Kantorovich-Rubenstein representation and [98, eq. (43)] generalises the set of transport plans  $\pi$  to complex measures, it is not entirely clear if a generalisation of the Wasserstein distance to complex measures  $\mu_1, \mu_2$

<sup>17</sup>If  $\mathcal{X}$  is a bounded metric space, we define its *diameter* by  $\text{diam}(\mathcal{X}) = \sup_{x, y \in \mathcal{X}} d(x, y)$ . The situation of two measures with equal total mass is often called *balanced* optimal transport.

has been considered before. However,  $p$ -Wasserstein distances for non-normalised, signed measures have been introduced by modification of the primal problem in [130, 131]. Additionally, an approach similar to ours starting from the Kantorovich-Rubenstein formula is studied in [96, Def. 8] and [97]. There, the unbalanced optimal transport distance between signed measures  $\mu, \nu$  with finite total variation is defined as

$$\sup_{\substack{f: \text{Lip}(f) \leq 1 \\ \|f\|_\infty \leq \tau}} \int_{\mathcal{X}} f(x)(d\mu(x) - d\nu(x)) \quad (1.10)$$

for some  $\tau > 0$ , see [96, Def. 8]. For  $\tau \rightarrow \infty$  and two measures with equal mass this directly yields the classical 1-Wasserstein distance. On the other hand, the supremum norm of an optimiser for (1.10) can be bounded in terms of the diameter of  $\mathcal{X}$  if the set  $\mathcal{X}$  is bounded and  $\mu$  and  $\nu$  have equal mass since the condition on the Lipschitz constant bounds the possible growth of the function. Therefore, we can extend the Wasserstein distance as follows.

**Proposition 1.4.5** (Complex, unbalanced 1-Wasserstein distance). *Let  $\mathcal{X}$  be compact. Then, the function  $W_1 : \mathcal{M}(\mathcal{X}) \times \mathcal{M}(\mathcal{X}) \rightarrow \mathbb{R}_{\geq 0}$ ,*

$$W_1(\mu, \nu) := \sup_{\substack{f: \text{Lip}(f) \leq 1 \\ \|f\|_\infty \leq \frac{1}{2} \text{diam}(\mathcal{X})}} \left| \int_{\mathcal{X}} f(x)(d\mu(x) - d\nu(x)) \right|,$$

*defines a metric called complex, unbalanced 1-Wasserstein distance on  $\mathcal{M}(\mathcal{X})$  that agrees with the definition of a balanced 1-Wasserstein distance from Definition 1.4.3 on  $\mathcal{M}_{+,1}(\mathcal{X})$ .*

*Proof.* Nonnegativity, symmetry and the triangle inequality are trivial. If one considers probability measures  $\mu, \nu \in \mathcal{M}_{+,1}(\mathcal{X})$ , we can always add a constant  $c \in \mathbb{R}$  to  $f$  without changing the value of  $\int_{\mathcal{X}} f(x)(d\mu(x) - d\nu(x))$ . The real valued, continuous function  $f$  attains its minimum  $f_{\min}$  and maximum  $f_{\max}$  on the compact set  $\mathcal{X}$  such that

$$\begin{aligned} \left| f(x) - \frac{f_{\min} + f_{\max}}{2} \right| &\leq \max \left( \left| f_{\min} - \frac{f_{\min} + f_{\max}}{2} \right|, \left| f_{\max} - \frac{f_{\min} + f_{\max}}{2} \right| \right) \\ &\leq \frac{1}{2} |f_{\min} - f_{\max}| \leq \frac{\text{diam}(\mathcal{X})}{2}. \end{aligned}$$

Thus, the condition  $\|f\|_\infty \leq \frac{\text{diam}(\mathcal{X})}{2}$  can be neglected for  $\mu, \nu \in \mathcal{M}_{+,1}(\mathcal{X})$  and the proposed extension of  $W_1$  agrees with Definition 1.4.3 on  $\mathcal{M}_{+,1}(\mathcal{X})$ .

In order to show that  $W_1$  is definite, set  $\tilde{\mu} = \mu - \nu$  and let  $\int_{\mathcal{X}} f(x) d\tilde{\mu}(x) = 0$  for all Lipschitz continuous functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  and assume that there is a measurable set  $A \subset \mathcal{X}$  with  $\tilde{\mu}(A) \neq 0$ . Since  $\tilde{\mu}$  is a Borel measure, we can consider  $A$  being closed. Denoting the projection operator on  $A$  by  $\text{proj}_A(x) := \text{argmin}_{y \in A} \|x - y\|_{\mathcal{X}}$ , we define for any  $\epsilon > 0$  a function  $f_\epsilon : \mathcal{X} \rightarrow \mathbb{R}$ ,

$$f_\epsilon(x) = \begin{cases} 1, & x \in A, \\ 1 - \frac{1}{\epsilon} \|\text{proj}_A(x) - x\|_{\mathcal{X}}, & x \in ((A + B_\epsilon(0)) \setminus A) \cap \mathcal{X}, \\ 0, & \text{else.} \end{cases}$$

It is straightforward to show that  $\text{Lip}(f_\epsilon) \leq \epsilon^{-1}$  and therefore  $0 = \int_{\mathcal{X}} f_\epsilon d\tilde{\mu}$  leading to the contradiction  $\mu(A) = 0$  by applying the dominated convergence theorem.  $\square$

## 1 Preliminaries

Proposition 1.4.5 can immediately be applied to  $\mathcal{X} = \mathbb{T}^d$  with the *wrap around metric*

$$\|t_1 - t_2\|_{\mathbb{T}^d} = \min_{j \in \mathbb{Z}^d} \|t_1 - t_2 + j\|_2$$

and we will use the resulting complex, unbalanced 1-Wasserstein distance for  $\mathcal{M}(\mathbb{T}^d)$ . Therein, the diameter of  $\mathcal{X} = \mathbb{T}^d$  is  $\text{diam}(\mathbb{T}^d) = \frac{\sqrt{d}}{2}$  and this leads to

$$W_1(\mu, \nu) := \sup_{\substack{f: \text{Lip}(f) \leq 1 \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}}} \left| \int_{\mathcal{X}} f(x) (\mathrm{d}\mu(x) - \mathrm{d}\nu(x)) \right|$$

for  $\mu, \nu \in \mathcal{M}(\mathbb{T}^d)$ .

## 2 Condition of sparse super resolution

In many parts, this chapter is a summary and extension of the papers [67, 68]. In particular, Subsection 2.2.4 goes beyond the previously published work of the author.

As already mentioned in the introduction, super resolution is about the computation of a measure from its noisy, low pass version or equivalently from its perturbed Fourier coefficients. Instead of analysing the performance of specific algorithms in the case of noise, we are interested in the nature of this inverse problem itself, i.e. we want to study how badly even the “best” algorithm in some sense amplifies errors in the data. In this work, we want to restrict ourselves to the case of a sparse ground truth measure consisting of a linear combination of Dirac measures.<sup>18</sup> Then, the task is to recover an underlying measure  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  with a *node set*  $Y \subset [0, 1]^d$ ,  $d \in \mathbb{N}_{\geq 1}$ , and *weights*  $(\alpha_t)_{t \in Y} \in \mathbb{C}^{|Y|}$  from noisy samples of its convolution with the *point spread function* (PSF) of the optical system  $h \in L^1(\mathbb{R}^d) \cap C_0(\mathbb{R}^d)$ .<sup>19,20</sup> In other words, one has access to perturbed values of

$$g(x) = (h * \mu)(x) = \sum_{t \in Y} \alpha_t h(x - t) \quad (2.1)$$

for  $x$  evaluated on some pixel grid in  $[0, 1]^d$ . By computing or more practically approximating the Fourier transform using the samples, one obtains estimates for the Fourier coefficients or moments of the measure given by

$$\hat{\mu}(k) = \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k}$$

for discrete frequencies  $k \in \mathbb{Z}^d$  such that  $\hat{h}(k) \neq 0$ . More realistically, the set of trustworthy spectral information would contain all  $k \in \mathbb{Z}^d$  where  $\hat{h}(k)$  is large compared to the noise level. As it is a widely used assumption to consider a radial PSF  $h$  and likewise a radial *optical transfer function* (OTF)  $\hat{h}$  (e.g. cf. [75, 28]), we assume access to estimates

$$\hat{\tilde{\mu}}(k) = \hat{\mu}(k) + \hat{\rho}(k) = \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k} + \hat{\rho}(k), \quad k \in \mathbb{Z}^d, \|k\|_2 \leq n, \quad (2.2)$$

for some  $n > 0$  and deterministic noise  $|\hat{\rho}(k)| \leq \varrho$  with noise level  $\varrho > 0$  on the Fourier coefficients. In Subsection 2.2.3 and Subsection 2.2.4, we will later assume also random noise on the moments or deterministic noise in the spatial domain.

Now, we want to study how much this noise is amplified if we compare the noise level in the data to the distance of the recovered parameters of the measure to the original one. Therefore, we introduce a notion for distance in the data and the parameter space.

<sup>18</sup>This assumption does not only include many single molecule microscopy applications but sparse measures might also be good approximations to measures with a more complicated support (e.g. cf. Section 4.1).

<sup>19</sup>For a set  $Y$  we denote its cardinality by  $|Y|$ .

<sup>20</sup>We have seen the bandlimited Airy function as a prototypical PSF in the introduction. Another popular choice would be an appropriately chosen Gaussian function (e.g. cf. [122, 96]). Both types of choice fulfil  $h \in L^1(\mathbb{R}^d) \cap C_0(\mathbb{R}^d)$ .

## 2 Condition of sparse super resolution

We call the measure  $M$ -sparse if  $\mu$  has up to  $M$  nodes and we interpret the nodes as elements of the  $d$ -dimensional torus  $\mathbb{T}^d$  since (2.2) allows for a 1-periodic ambiguity in the nodes. Therefore, the distance between nodes should be the *wrap around distance*  $\|t_1 - t_2\|_{\mathbb{T}^d} = \min_{j \in \mathbb{Z}^d} \|t_1 - t_2 + j\|_2$  and the Euclidean *wrap-around separation* of some finite set  $Y$  is then defined as

$$\text{sep } Y := \min_{t_1, t_2 \in Y, t_1 \neq t_2} \|t_1 - t_2\|_{\mathbb{T}^d} = \min_{t_1, t_2 \in Y, t_1 \neq t_2} \min_{j \in \mathbb{Z}^d} \|t_1 - t_2 + j\|_2.$$

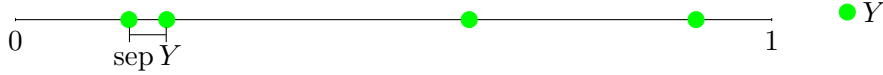


Figure 2.1: Definition of separation distance on  $\mathbb{T}$  for some finite set  $Y$ .

For  $q > 0$ , we define the *set of  $q$ -separated complex measures*

$$\mathcal{M}(q) := \left\{ \sum_{t \in Y} \alpha_t \delta_t : \alpha_t \in \mathbb{C}, Y \subset \mathbb{T}^d \text{ finite, sep } Y \geq q \right\}$$

as a subset of the space  $\mathcal{M}(\mathbb{T}^d)$ . The corresponding *truncated moment set* or *set of exponential sums* is

$$\widehat{\mathcal{M}}^n(q) := \left\{ \left( \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k} \right)_{k \in \mathbb{Z}^d: \|k\|_2 \leq n} : \sum_{t \in Y} \alpha_t \delta_t \in \mathcal{M}(q) \right\},$$

where we again restrict ourselves to the sampling set  $\{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}$ . The standard Euclidean norm  $\|\cdot\|_2$ , where  $\|\hat{\mu}\|_2^2 := \sum_{k \in \mathbb{Z}^d: \|k\|_2 \leq n} |\hat{\mu}(k)|^2$ , always induces a metric on  $\widehat{\mathcal{M}}^n(q)$  and this manifold is sketched in Figure 2.2. The shape of  $\widehat{\mathcal{M}}^n(q)$  determines how large the ratio of the noise  $\hat{\rho}$  and the distance between ground truth  $\mu$  and the best sparse

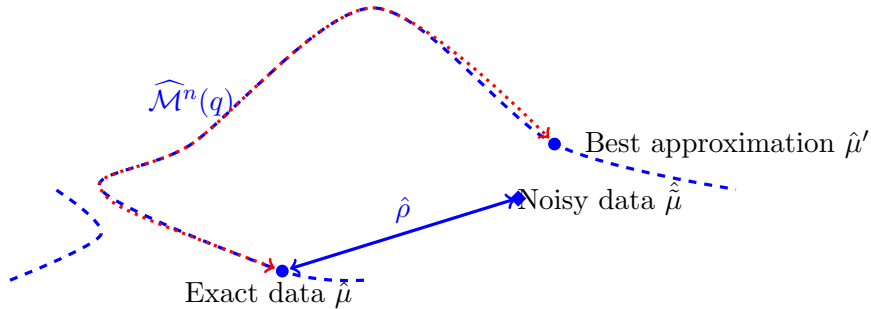


Figure 2.2: Sketch of the manifold  $\widehat{\mathcal{M}}^n(q)$  and idea of noise amplification as the worst ratio between the difference of  $\mu$  and  $\mu'$  measured by difference of their parameters (red arrow) compared to noise level  $\hat{\rho}$  in  $\widehat{\mathcal{M}}^n(q)$  (blue arrow). As the condition on the separation is usually chosen such that  $|Y|$  is noticeably smaller than the number of moments, the manifold is typically a proper subset of  $\mathbb{C}^{\{k \in \mathbb{Z}^d: \|k\|_2 \leq n\}}$ .



approximation  $\mu'$  might be.<sup>21</sup> We study this question in this chapter by summarising related results at first. Afterwards, we develop a framework in which the condition number of the nonlinear mapping of Fourier coefficients to the closest measure is analysed. This includes the construction of a new *minorant function* in Subsection 2.2.1 which is used to prove an inequality for the difference of nodes and weights in Subsection 2.2.2. Based on this, we study the diffraction limit as a transition of the condition number from polynomial to exponential with respect to the bandlimit  $n$  in Subsection 2.2.3. Our findings can then be applied to the condition of the full inverse problem starting with image data instead of Fourier coefficients in Subsection 2.2.4. Finally, we conclude a result for the smallest singular value of Vandermonde matrices with pair clusters in Section 2.3. This finding comes with the advantage that it is applicable in a multivariate setting if the clusters are separated by a term independent of the distance of the clustering pairs.

## 2.1 Summary of related results

As mentioned in Section 1.1, one has to distinguish the stability of an algorithm from the condition of the underlying problem. Nevertheless, the stability of an algorithm can not be better than the condition of the task such that results about stability of a particular algorithm already yield an upper bound for the condition. For the sparse super resolution problem, there exist a wide range of algorithms which can be divided into two types of algorithms. On one hand, there are variational methods which analyse the problem by finding a solution of certain optimisation schemes. For example, Candès and Fernandez-Granda [52, 22] promoted the use of a  $\ell^1$ -regularised minimisation problem yielding sparse reconstructions and established theory about it saying that this approach stably recovers the ground truth measure in the univariate setting if  $qn \geq 1.26$  ([52]). Similarly, the work by Duval and Peyré [43] has shown robust reconstruction with the Beurling-Lasso (BLASSO) under the hypothesis of a “non-degenerate source condition” being itself linked to an assumption on the minimal separation in relation to the cutoff frequency  $n$ , see [43, p. 1331].<sup>22</sup> The same authors together with Denoyelle [34] additionally obtained stability independent of the separation distance under the prior information that one is confronted with positive measures. More recently, there were also attempts to use variational methods of deep learning but for methods like *Deep STORM*, see [122, 121], (almost) no theoretical guarantees for the stability are known yet.

On the other hand, parametric methods commonly summarised by the term *subspace methods* are also used for the super resolution problem and there exists a long list of attempts to analyse stability of subspace methods including Li et al. [99], Liao et al. [101], Aubel and Bölcskei [7], Potts und Tasche [135] and Sahnoun [140]. As an example, we briefly describe the results of Fan and Li, cf. [49]. They present a variant of Prony’s method that approximates the ground truth measure with an error proportional to  $Md\rho^{\mathcal{O}(1/M)}$  in the Wasserstein distance with high probability. Here,  $\rho$  is the level of noise,  $M$  the number of nodes satisfying  $n = 2M - 1$  and  $d$  the dimension. This result can be seen as

<sup>21</sup>In particular, we will find by Theorem 3.1.5 that  $\widehat{\mathcal{M}}^n(q)$  for  $q$  sufficiently large is not path-connected in general because this theorem allows to conclude that elements in  $\widehat{\mathcal{M}}^n(q)$  which are very close to each other have the same number of parameters  $(t, \alpha_t)_{t \in Y}$ . Therefore, elements of  $\widehat{\mathcal{M}}^n(q)$  can only be connected by a path if they have the same number of parameters such that different connected components must exist. This non-connectedness is indicated in Figure 2.2.

<sup>22</sup>Here, the Beurling-Lasso describes a  $\ell^1$ -regularised optimisation scheme which is formulated on the space of Radon measures and thus independent of any choice of discretisation of the domain.

## 2 Condition of sparse super resolution

a stability estimate for a particular algorithm in the strong super resolution case as there is no separation condition involved and the rate deteriorates with an increasing number of nodes  $M$ . The latter was already observed by Donoho, see [39, Thm. 1.3], and can also be observed for variational methods, see [43].

Another way to estimate the condition of sparse super resolution is given through the analysis of the smallest singular value of Vandermonde matrices which was done among others by Moitra [117], Aubel and Bölcskei [6], Nagel et al. [88, 120], Batenkov et al. [10], and the references therein. The ratio behind this observation is that a simplification of the super resolution problem would be that the nodes are already known. In that case, it is natural to set up a Vandermonde matrix based on the given nodes and to recover the weights by a least square approach. Then, the noise amplification is governed by the condition number of the Vandermonde matrix and the latter depends heavily on the size of the smallest singular value. If nodes of the Vandermonde matrix are very close, the smallest singular value is very small whereas the condition number can be nicely bounded if the node set is *well separated*. Thus this qualitative transition gives already an idea of a diffraction limit distinguishing between strong and weak super resolution. Moreover, one can use these estimates in order to study the stability of the node recovery using algorithms like ESPRIT or matrix pencil, see [7, 120]. Nevertheless, considering the recovery of nodes and weights separately is a simplification of the problem and the condition of a problem is not the same as the stability of an algorithm.<sup>23</sup> Our approach to analyse the condition of the complete task is more general and allows to conclude a result for the smallest singular values of Vandermonde matrices with pairwise clustering nodes, see Section 2.3.

An interesting statistical approach was presented by Ferreira Da Costa et al. [53] based on the *Cramér-Rao (CR) lower bound*. There, the CR lower bound means the inequality that the covariance matrix of any unbiased estimator  $\hat{\theta}$  for a vector of parameters  $\theta$  can be bounded from below by the inverse of the so called *Fisher information matrix*  $J(\theta)$ , in other words we have

$$\mathbb{E} \left[ (\hat{\theta}(y) - \theta)(\hat{\theta}(y) - \theta)^\top \right] \succeq J(\theta)^{-1}$$

where  $A \succeq B$  for Hermitian matrices  $A, B \in \mathbb{C}^{m \times m}$  is meant in the sense that  $A - B$  is positive semidefinite, e.g. cf. [126]. In the context of super resolution, it is obvious that  $\theta$  consists of all nodes and weights and the Fisher information matrix calculated for the univariate problem in [53] has a structure similar to a *confluent Vandermonde matrix* defined in [57]. Hence, lower bounds on the covariance of any unbiased estimator for the parameters of the sparse measure can be derived through the smallest singular value of the confluent Vandermonde matrices and the latter is done for the univariate setting in [53] by using minorant and majorant functions derived from the Beurling-Selberg extremal functions. This results in a theorem, see [53, Prop. 3], stating that the one-dimensional super resolution problem is stable in the sense that the Fisher information matrix is well conditioned if the separation  $q$  of the nodes satisfies  $(2n + 1) \cdot q > 3.54$ . It is conjectured that the optimal lower bound for  $(2n + 1) \cdot q$  is two and we will prove this in Subsection 2.2.3 together with a multivariate extension. Furthermore, this generalisation to higher dimensions will be consistent with our notions for stability and the diffraction limit.

<sup>23</sup>One might argue to close the gap between condition and stability by the Cramér-Rao lower bound discussed in the next paragraph. This bound allows to control the minimal possible variance of an estimator and can therefore be able to show optimal performance of an algorithm such that understanding the stability of the best algorithm already explains the condition of the underlying problem.

A different approach discussing the correct identification of the number and position of nodes was given by Liu et al. in [109, 103]. Working in the case  $d = 1$  (cf. [109]) and  $d \geq 2$  (cf. [108]) with an interval or a ball of sufficient size instead of  $\mathbb{T}^d$ , two different kinds of resolution limits are defined. On one hand, the *computational resolution limit to the number detection problem* is considered as the minimal separation  $q_{\text{Liu, number}}$  depending on cut-off-frequency  $n$ , smallest weight  $\alpha_{\min}$ , noise level  $\varrho$  and overall number of nodes  $M$  such that for any  $M$ -sparse ground truth measure  $\mu_0 \in \mathcal{M}(q_{\text{Liu, number}})$  and for every perturbation by noise of size  $\varrho$  the resulting noisy Fourier coefficients of the ground truth have the property that every sparse measure agreeing with the measured Fourier coefficients up to this error  $\varrho$  has at least  $M$  nodes, see [108, Def. 2.2].<sup>24</sup> In other words, this resolution limit guarantees that every reconstruction, which is compatible with the noise, consists of the correct number of parameters or more. This notion of a resolution limit can then be estimated by

$$\frac{1}{n} \left[ \frac{\varrho}{\alpha_{\min}} \right]^{1/(2M-2)} \lesssim q_{\text{Liu, number}} \lesssim_d \frac{M}{n} \left[ \frac{\varrho}{\alpha_{\min}} \right]^{1/(2M-2)}, \quad (2.3)$$

see [108, Thm. 2.3 and Pro. 2.4].<sup>25</sup> On the other hand, another resolution limit is suggested in the same work by demanding a sparse recovery where each recovered node is close to exactly one ground truth node. More formally, the *computational resolution limit to the support recovery problem* is defined as the minimal separation  $q_{\text{Liu, support}}$  such that for any  $M$ -sparse ground truth measure  $\mu_0 \in \mathcal{M}(q_{\text{Liu, support}})$  and for every perturbation by noise of size  $\varrho$  the resulting noisy Fourier coefficients of the ground truth have the property that there exists a neighbourhood size  $\delta'$  such that every  $M$ -sparse measure agreeing with the measured Fourier coefficients up to error  $\varrho$  consists of nodes having exactly one ground truth node in their neighbourhood of radius  $\delta'$ , see [108, Def. 2.6]. In [108, Thm. 2.7 and Pro. 2.8], this is then bounded as

$$\frac{1}{n} \left[ \frac{\delta}{\alpha_{\min}} \right]^{1/(2M-1)} \lesssim q_{\text{Liu, support}} \lesssim_d \frac{M}{n} \left[ \frac{\delta}{\alpha_{\min}} \right]^{1/(2M-1)} \quad (2.4)$$

where the upper bounds as shown in (2.3) and (2.4) are an improvement of Liu in [104]. If only positive measure were considered, the linear dependency of the upper bounds on  $M$  can be dropped and the orders match completely (cf. [106]). Furthermore, the influence of multiple illuminations can also be analysed in this framework and a gain in resolution is then obtained by incoherence of illumination matrix, see [107] and Section 4.1. As shown in [108] for the multivariate cases  $d = 2, 3$ , the two computational resolution limits introduced by Liu appear to govern the success or failure probability of parametric algorithms like the matrix pencil method and projected matrix pencil method. Moreover, the results give an insight into the condition of strong super resolution where the nodes might be very closely spaced. However, the geometry of the nodes is not considered and thus the exponent for the signal to noise ratio  $\frac{\alpha_{\min}}{\delta}$  deteriorates for a large number of nodes  $M$ . Additionally, the constants hid behind the “ $\lesssim$ ”-notation are not practically close to each other and thus it is difficult to apply the results to describe the actual resolution limit of an imaging system.

<sup>24</sup>The notions for distance and separation are similar to our setting even though this source does not work on the torus but on a compact domain of  $\mathbb{R}^d$ . Moreover, the weights of the sparse measures are real whereas our construction is for complex weights in general.

<sup>25</sup>We write  $a \lesssim b$  if  $ca \leq b$  for some constant  $c > 0$  and indicate dependency of  $c$  on some variable  $d$  by  $a \lesssim_d b$ .

## 2 Condition of sparse super resolution

In order to circumvent the dependency of stability on  $M$  for strong super resolution, information on the geometry of the nodes needs to be taken into account. More precisely, one can see that the exponent  $M$  can be replaced by a smaller  $\ell$  representing the largest cluster size in the node set.<sup>26</sup> In the univariate situation and for continuous measurements, this was taken into account by the work of Batenkov et al. [12]. Recently, a specialisation by Liu and Ammari in [105] to positive instead of complex measures was made. Since the results are similar, we focus on [12]. The paper analyses the minmax error for the recovery of parameters of sparse measures from noisy data depending not only on noise level  $\varrho$  and frequency-sampling cutoff  $n$  but also on the question whether the particular node is contained in a cluster of  $\ell > 1$  nodes. More precisely, the best recovery algorithm recovers the non-clustered nodes and corresponding weights with error  $\mathcal{O}(\frac{\varrho}{n})$  or  $\mathcal{O}(\varrho)$  respectively whereas the nodes in a cluster with  $\ell$  elements are computed with error  $\mathcal{O}(\frac{\varrho}{n}(nq)^{-2\ell+2})$  where  $q < n^{-1}$  is the minimal separation of the sparse ground truth. Additionally, the error for the recovery of weights corresponding to clustering nodes admits the order  $\mathcal{O}(\varrho(nq)^{-2\ell+1})$ , cf. [12, Thm. 2.8]. The proof uses some sort of a “quantitative inverse function theorem” (cf. [12, Thm. B.1]) and seems to be limited to the one-dimensional case. This drawback together with their remark that “it is of great interest” to get control over the “problem condition number” [12, p. 2] was one of the main motivations for this work since the term “problem condition number” was not explicitly defined in [12].

On the theoretical side, our methods are based on the work by Diederichs [37, 38] where minorant functions in  $d = 1, 2$  are established in order to bound the difference in Fourier data from below by differences in the parameters. Using minorants being maximal in the origin in the spatial and in Fourier domain, Diederichs was able to derive the same orders for the recovery error as [12] in  $d = 1$  for well separated nodes and can extend them to  $d = 2$ . We improve this by introducing minorant functions with optimal localisation yielding a generalisation to arbitrary dimensions.

Another important contribution to the analysis of the condition of super resolution was given by Chen and Moitra where the resolution limit was described as a transition of the “sample complexity” for recovering sparse measures “from polynomial to exponential” (cf. [28, p. 3]). They prove that this transition in the condition of the two-dimensional problem happens at a separation distance  $q$  with  $1.15 \approx \sqrt{\frac{4}{3}} \leq q \cdot n \leq \frac{2j_{0,1}}{\pi} \approx 1.53$ . The used method is again driven by a minorant function being radial in their analysis. Utilising a different radial function, we are able to improve the upper bound for the location of the transition point to approximately 1.22.

In contrast to the previous works describing super resolution as the recovery of a measure from its low pass Fourier coefficients, the works of Eftekhari et al. [45, 44] for  $d = 1$  or  $d = 2$  respectively and its multivariate extension by Kurmanbek et al. [90] consider the recovery of a measure from its convolution with some point spread function (PSF). More precisely, all papers assume access to (noisy) convolution measurements  $g$  with PSF  $h$ , i.e.  $\|h * \mu - g\| \leq \varrho$  for some noise level  $\varrho > 0$  and a nonnegative ground truth measure  $\mu \in \mathcal{M}_+([0, 1]^d)$ . This assumption motivates to study the convex feasibility problem of finding

$$\mu' \in \mathcal{M}_+([0, 1]^d) \text{ s.t. } \|h * \mu' - g\| \leq \varrho' \quad (2.5)$$

<sup>26</sup>Usually, nodes are said to form a cluster if they are separated by less than the Rayleigh limit meant in the sense that nodes belong to a cluster if they are separated by less than  $c_d n^{-1}$  where  $c_d > 0$  depends on the dimension  $d$ , e.g. see [12].

for some  $\varrho' > \varrho$ . Here, the considered norm is the 2-norm in univariate case [45] or the Frobenius norm for multivariate data sets ([44, 90]). The univariate analysis in [45] has a special focus on  $h$  being a Gaussian kernel and generalises this setting to systems that form a so called *Chebyshev system*. This is used in the multivariate extensions [44, 90] by considering a tensor product of Chebyshev kernels for the PSF. All papers do not focus on solving the feasibility problem (2.5) numerically but they use the generalised Wasserstein distance  $d_{GW}$  from [130, 131] in order to compare reconstructed measures with the ground truth where both measures not necessarily have the same mass.<sup>27</sup> The authors especially emphasise that (2.5) needs no sparsity constraint but only the assumption of nonnegativity. On one hand, all three papers consider the recovery of a sparse ground truth measure  $\mu$  with exact convolution measurements, i.e.  $\varrho = 0$ . Then, the feasibility problem for  $\varrho' = 0$  recovers  $\mu$  if  $h$  forms a Chebyshev system and the number of measurements in each dimension is larger than two times the sparsity of  $\mu$  ([45, Prop. 8], [44, Prop. 2] and [90, Thm. 2.2]). The proof bases on the construction of a dual certificate guaranteeing the success of (2.5). On the other hand, the more interesting case of positive noise level  $\varrho > 0$  and arbitrary ground truth measure  $\mu$  is for example studied as follows in [44, Thm. 11]. If  $\|h * (\mu_1 - \mu_2)\| \leq L d_{GW}(\mu_1, \mu_2)$ , e.g.  $L \approx \|\nabla h\|$ , and  $\varrho' \geq (1 + L \min_{\nu \in \mathcal{M}(q), |Y^\nu| \leq r} d_{GW}(\mu, \nu))\varrho$ , the solution to (2.5) recovers the ground truth up to an error

$$d_{GW}(\mu, \mu') \leq c_1 \varrho + c_2(q) + c_3 \min_{\nu \in \mathcal{M}(q), |Y^\nu| \leq M} d_{GW}(\mu, \nu), \quad (2.6)$$

where  $c_1, c_2, c_3$  are not completely explicit functions of  $q, M$  and the PSF  $h$ . Therefore,  $\mu'$  approximates  $\mu$  well if  $\mu$  can be well-approximated by a  $q$ -sparse measure and  $\varrho$  is small. However, one drawback of this result is the lack of an explicit presentation of  $c_2$ .<sup>28</sup> The authors just highlight  $c_2(0) = 0$  but certainly  $c_2(q) > 0$  for some  $q > 0$ . Hence, the estimate (2.6) does not allow to conclude that the reconstruction error for a sparse measure  $\mu$  with separation  $q > 0$  does go to zero as  $\varrho$  tends to zero because the second term in the upper bound remains. This behaviour might be seen as an additional disadvantage of this stability result for super resolution.

## 2.2 Condition estimates

### 2.2.1 Admissible functions

We want to use a function  $\psi$  with various properties to apply Poisson's summation formula in order to relate Fourier coefficients of a discrete measure  $\mu$  to its parameters in real space. As we need a minorant in Fourier domain, we are interested in functions  $\psi$  such that their Fourier transform  $\hat{\psi}$  is a minorant to the indicator function of the Euclidean unit ball. Together with various other assumptions, we call such a function *admissible*. Beyond the condition  $\psi(0) > 0$  used in [86], we additionally require similar to [37, 38] that this is the global maximum.

**Definition 2.2.1** (Admissible function). Let  $d \in \mathbb{N}$  and  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a function  $\psi \in L^1(\mathbb{R}^d)$  which

<sup>27</sup>It is easy to show that  $W_1(\mu, \nu) \leq \max(1, \text{diam}(\mathcal{X})/2) d_{GW}(\mu, \nu)$  for  $\mu, \nu \in \mathcal{M}_+(\mathcal{X})$ . However, an estimate in the other direction is not directly obvious and thus it not clear whether  $d_{GW}$  and our generalisation  $W_1$  are equivalent on  $\mathcal{M}_+(\mathcal{X})$ .

<sup>28</sup>It is explained that  $c_2$  depends on  $\|\mu\|_{TV}$  and on  $q$  by some function  $\alpha(q)$  and this function  $\alpha$  is then defined via a maximum over the dual certificate [44, p. 182 and p. 192].

## 2 Condition of sparse super resolution

- (i) is continuous with compact support, i.e.  $\psi \in C_c(\mathbb{R}^d)$ ,
- (ii) attains its global maximum  $\psi(0) > 0$  in the origin allowing to find  $c_d > 0$  such that the bound

$$\psi(0) - \psi(x) \geq c_d \|x\|_2^2$$

for any  $x \in \text{supp } \psi$  holds,

- (iii) and satisfies  $\hat{\psi}(v) \in \mathbb{R}$  for all  $v \in \mathbb{R}^d$  with sign

$$\hat{\psi}(v) \begin{cases} \geq 0 & \|v\|_2 \leq 1, \\ \leq 0 & \|v\|_2 \geq 1. \end{cases} \quad (2.7)$$

Then, we call a function  $\psi$  fulfilling (i)-(iii) *admissible*.

We have summarised classical results for functions satisfying (i) and (iii) in Section 1.3. Additionally, we included a univariate function from [37] that also meets condition (ii). Based on the idea from [83] explained already in (1.7), we find admissible functions in the general multivariate case. We mostly focus on the case where the support of an admissible function  $\psi$  is an Euclidean ball such that it is natural by symmetry to consider a radial function  $\psi$ .

**Lemma 2.2.2** (Support on a ball). *For  $d \geq 1$  we define  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ ,*

$$\varphi(x) = \begin{cases} 1 - \left( \frac{j_{d/2,1}}{2\pi \|x\|_2} \right)^{d/2-1} \frac{J_{d/2-1}(2\pi \|x\|_2)}{J_{d/2-1}(j_{d/2,1})}, & \|x\|_2 < \frac{j_{d/2,1}}{2\pi}, \\ 0, & \text{otherwise.} \end{cases} \quad (2.8)$$

Moreover, let  $\Delta = \sum_{s=1}^d \frac{\partial^2}{\partial x_s^2}$  be the Laplace operator and  $\psi_\tau: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ ,  $\tau \geq 0$ ,

$$\psi_\tau(x) = \left( \frac{1}{\sqrt{1+\tau}} \right)^d [4\pi^2(1+\tau) + \Delta] (\varphi * \varphi) \left( \frac{x}{\sqrt{1+\tau}} \right).$$

Then,  $\psi_\tau$  with  $\tau > 0$  is admissible, its support satisfies  $\text{supp } \psi_\tau = B_{q_\tau}(0)$  with

$$q_\tau := \sqrt{1+\tau} \frac{j_{d/2,1}}{\pi},$$

and there is a constant  $c_d > 0$  depending only on  $d$  that allows the estimate

$$\psi_\tau(0) - \psi_\tau(x) \geq c_d \tau (1+\tau)^{-d/2-1} \|x\|_2^2 \quad (2.9)$$

for all  $x \in \text{supp } \psi_\tau$ .

*Proof.* We directly find  $\text{supp } \psi_\tau = \{x \in \mathbb{R}^d : \|x\|_2 \leq \sqrt{1+\tau} \frac{j_{d/2,1}}{\pi}\}$  since  $\text{supp}(\varphi * \varphi) = \text{supp}(\varphi) + \text{supp}(\varphi)$  where the sum is interpreted as a Minkowski sum. Moreover,  $\varphi$  is continuous by construction and we note that the window function  $\varphi$  in (2.8) admits

$$(4\pi^2 + \Delta)\varphi = \gamma \cdot \mathbb{1}_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} \quad (2.10)$$

for some constant  $\gamma > 0$ , see [29]. Hence,  $\psi_0$  as a convolution of the integrable function  $\mathbb{1}_{B_{q_0}(0)}$  with the continuous function  $\varphi$  is again continuous and thus the same holds for  $\psi_\tau$ .

Therefore, condition (i) for an admissible function is fulfilled. Using Proposition 1.2.2, the Fourier transform  $\hat{\psi}$  can be expressed as

$$\hat{\psi}_\tau(v) = 4\pi^2(1+\tau) [1 - \|v\|_2^2] [\hat{\varphi}(\sqrt{1+\tau}v)]^2 \quad (2.11)$$

and this gives (2.7) or condition (iii) respectively. Furthermore, one can directly deduce  $J_{d/2-1}(j_{d/2,1}) < 0$ ,  $\varphi \geq 0$ , and  $\frac{\partial\varphi}{\partial r} \leq 0$  by Lemma 1.3.2 such that one observes

$$\frac{\partial\psi_0}{\partial r} = \gamma \cdot \mathbb{1}_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} * \frac{\partial\varphi}{\partial r} \leq 0$$

and thus  $\psi_0(x)$  has a maximum at  $x = 0$ . Finally, Hölder's inequality gives

$$(\varphi * \varphi)(t) < \left( \int_{\mathbb{R}^d} \varphi(s)^2 ds \right)^{1/2} \left( \int_{\mathbb{R}^d} \varphi(t-s)^2 ds \right)^{1/2} = (\varphi * \varphi)(0)$$

and the inequality is strict because  $\varphi(s)$  and  $\varphi(t-s)$  are not constant multiples of each other if  $t \neq 0$ . This means already that  $\psi_\tau$  is maximal in zero. For the rate, we use the inverse Fourier transform of the radial function  $\hat{\varphi}$  as presented in (1.1), i.e.,

$$(\varphi * \varphi)(0) - (\varphi * \varphi)(x) = \int_0^\infty \frac{|\hat{\varphi}(\omega e_1)|^2 \omega^{d/2}}{(2\pi)^{-1}} \left\{ \frac{(\pi\omega)^{d/2-1}}{\Gamma(\frac{d}{2})} - \frac{J_{d/2-1}(2\pi\|x\|_2\omega)}{\|x\|_2^{d/2-1}} \right\} d\omega, \quad (2.12)$$

where  $e_1 = (1, 0, \dots, 0) \in \mathbb{R}^d$  denotes the first unit vector. By applying the Fourier transform on both sides of (2.10), we see that  $\hat{\varphi}(\omega e_1)$  is proportional to  $\frac{J_{d/2}(j_{d/2,1}\omega)}{\omega^{d/2}(1-\omega^2)}$ .<sup>29</sup> Consequently, the function  $\hat{\varphi}$  has sufficient decay such that we can conclude  $\mathcal{O}(\omega^{-6})$  as the order of the integrand in (2.12).<sup>30</sup> Then, the dominated convergence theorem and the uniform convergence of the series defining the Bessel function allow for

$$\lim_{\|x\|_2 \rightarrow 0} \frac{(\varphi * \varphi)(0) - (\varphi * \varphi)(x)}{\|x\|_2^2} = 2\pi \int_0^\infty |\hat{\varphi}(\omega e_1)|^2 \frac{\pi^{d/2+1} \omega^{d+1}}{\Gamma(\frac{d}{2} + 1)} d\omega > 0.$$

By the regularity of  $\varphi * \varphi$ , there exist  $c_0, r_0 > 0$  such that  $(\varphi * \varphi)(0) - (\varphi * \varphi)(x) \geq \frac{c_0}{4} \|x\|_2^2$  for  $\|x\|_2 \leq r_0$ . For  $r_0 \leq \|x\|_2 \leq \frac{j_{d/2,1}}{\pi}$  the function  $\varphi * \varphi$  has a maximal value  $c'_0 < (\varphi * \varphi)(0)$ . Summing up, we then end with

$$\begin{aligned} \psi_\tau(0) - \psi_\tau(x) &= \psi_0(0) - \psi_0(x) + 4\pi^2\tau \left( \frac{1}{\sqrt{1+\tau}} \right)^d \left[ (\varphi * \varphi)(0) - (\varphi * \varphi)((1+\tau)^{-1/2}x) \right] \\ &\geq 4\pi^2\tau \left( \frac{1}{\sqrt{1+\tau}} \right)^d \min \left( \frac{c_0}{4(1+\tau)} \|x\|_2^2, (\varphi * \varphi)(0) - c'_0 \right) \\ &\geq 4\pi^2\tau \left( \frac{1}{\sqrt{1+\tau}} \right)^{d+2} \min \left( \frac{c_0}{4}, \frac{(\varphi * \varphi)(0) - c'_0}{\left( \frac{j_{d/2,1}}{2\pi} \right)^2} \right) \|x\|_2^2 \end{aligned}$$

and setting  $c_d := 4\pi^2 \min \left( \frac{c_0}{4}, [(\varphi * \varphi)(0) - c'_0] \frac{4\pi^2}{j_{d/2,1}^2} \right)$  finishes the proof.  $\square$

<sup>29</sup>The Fourier transform of the indicator function of a ball is well-known and dividing by  $1-\omega^2$  originating from the differential operator  $4\pi^2 + \Delta$  needs to be understood in the limit sense as  $\omega \rightarrow 1$ .

<sup>30</sup>By the asymptotic expansion from Lemma 1.3.2,  $\hat{\varphi}$  has the order  $\omega^{-d/2-\frac{5}{2}}$  as  $\omega \rightarrow \infty$ . This implies  $\hat{\varphi}(\omega)^2 \in \mathcal{O}(\omega^{-d-5})$  yielding  $\varphi * \varphi \in C^4(\mathbb{R}^d)$  by [61, Prop. 3.3.12 or Ex. 2.4.1]. From (2.11), we observe  $\hat{\psi}_\tau \in \mathcal{O}(\omega^{-d-3})$  and conclude  $\psi_\tau \in C^2(\mathbb{R}^d)$  analogously.

## 2 Condition of sparse super resolution

We remark that  $\varphi$  as in (2.8) was proposed in [29, 60] as a building block for  $\psi_\tau$  with  $\tau = 0$ . However,  $\psi_0$  does not allow an estimate of the form (2.9) because one can see that the second order terms vanish as  $\tau \rightarrow 0$ . In fact, one can show for some constant  $\tilde{c}_d$  that

$$\lim_{x \rightarrow 0} \frac{\psi_0(0) - \psi_0(x)}{\|x\|_2^2} = \int_0^\infty (1 - \omega^2) \hat{\varphi}(\omega)^2 \frac{8\pi^{d/2+4}\omega^{d+1}}{\Gamma(d/2 + 1)} d\omega = \tilde{c}_d \int_0^\infty \frac{\omega J_{d/2}^2(j_{d/2,1}\omega)}{1 - \omega^2} d\omega = 0$$

where the last equality is an elegant and direct consequence of [154, p. 429, eq. (3)]. Therefore  $\psi_0$  does not admit a quadratic lower bound for  $\psi_0(0) - \psi_0(x)$ . We see in the following remark that a support with radius larger than the critical radius  $\frac{j_{d/2,1}}{\pi}$  is indeed necessary for this property.

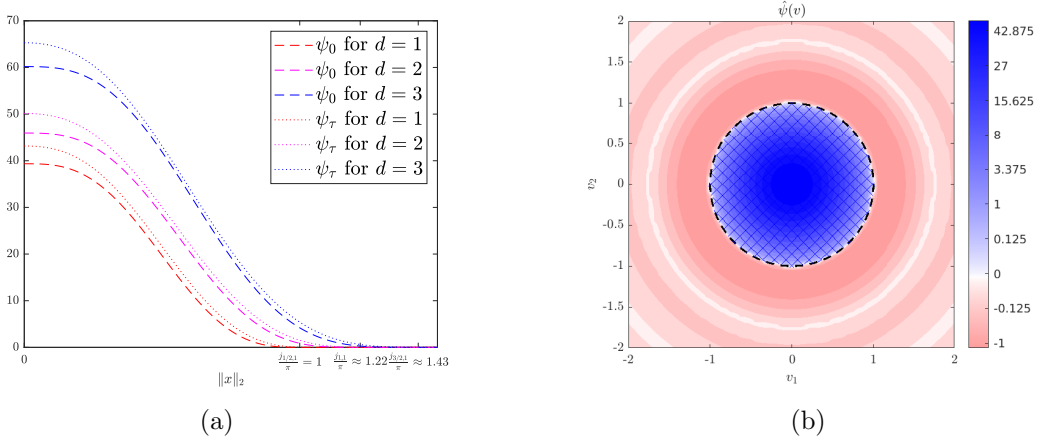


Figure 2.3: Admissible function  $\psi_\tau(x)$  for  $\tau = 0.1$  and non-admissible  $\psi_0(x)$  for  $d = 1, 2, 3$  as a function of  $\|x\|_2$  (a). One can at least imagine from its graph that the second derivative of  $\psi_0$  at  $x = 0$  vanishes such that it is not admissible according to our definition. Additionally, we display  $\hat{\psi}_{0.1}$  for  $d = 2$  and highlight its radial dependency as well as the change of the sign at  $\|v\|_2 = 1$  (b).

**Remark 2.2.3** (Optimality of the support). If a function  $\psi \in C^2(\mathbb{R}^d)$  is admissible and  $\text{supp } \psi \subset B_q(0)$ , then the quadratic estimate (2.9) implies that the Hessian of  $\psi$  satisfies  $\text{Hess } \psi(0) \prec 0$  and thus we have by Proposition 1.2.2

$$-4\pi^2 \int_{\mathbb{R}^d} \|v\|_2^2 \hat{\psi}(v) dv = [\Delta \psi](0) = \text{Tr}(\text{Hess } \psi(0)) < 0 \quad (2.13)$$

where  $\text{Tr}(A)$  denotes the *trace* of a matrix  $A$ . We can restrict to radial functions  $\psi$ , see [71, Lemmas 18 and 19], for which the latter condition (2.13) implies  $q > \frac{j_{d/2,1}}{\pi}$ , see [59, Thm. 1]. Lemma 2.2.2 shows that this critical separation is also sufficient and describes an admissible function with optimally small support in an explicit form.

For  $d \in \{1, 2, 3\}$ , we display the radial dependency of  $\psi_0$  and  $\psi_\tau$  for  $\tau = 0.1$  in Figure 2.3 and note in passing that  $\frac{j_{d/2,1}}{\pi} \approx 1, 1.22, 1.43$  for  $d = 1, 2, 3$ . Beyond being admissible, the function  $\psi_\tau$  has the remarkable property that even if we evaluate  $\psi$  at any finite number of points such that these points are separated by at least  $q_\tau$ , the sum of the evaluations of  $\psi_\tau$  is always smaller than the global maximum of  $\psi_\tau$ . Moreover, the difference of the global maximum and the sum of evaluations can even be bounded from below in terms of the square of the smallest radius of any of the finite sampling points. We formulate this in the following lemma.



**Lemma 2.2.4.** *Let  $\psi_\tau$  and  $q_\tau$  be defined as in Lemma 2.2.2. Then, there exists a constant  $c'_d > 0$  such that for finite  $Y \subset \mathbb{R}^d$  with  $\text{sep } Y := \min_{t,s \in Y, t \neq s} \|t - s\|_2 \geq q_\tau$  and for any  $t \in \mathbb{R}^d$  the inequality*

$$\psi_\tau(0) - \sum_{t' \in Y} \psi_\tau(t - t') \geq \begin{cases} c'_d \tau (1 + \tau)^{-d/2-1} \text{dist}(t, Y)^2, & \text{dist}(t, Y) \leq q_\tau, \\ \psi_\tau(0), & \text{dist}(t, Y) \geq q_\tau, \end{cases} \quad (2.14)$$

holds. Here,  $\text{dist}(t, Y)$  denotes the distance  $\text{dist}(t, Y) := \min_{s \in Y} \|t - s\|_2$ .

*Proof.* Due to the compact support of  $\psi_\tau$ , the inequality (2.14) is trivially fulfilled for  $t$  with  $\text{dist}(t, Y) = 0$  or  $\text{dist}(t, Y) \geq q_\tau$ . Therefore, let  $t \in \mathbb{R}^d$  be such that  $0 < \text{dist}(t, Y) < q_\tau$  and define the set  $Y_t = \{s = (1 + \tau)^{-1/2}(t - t'), t' \in Y\}$ . We directly find  $\text{sep } Y_t \geq \frac{j_{d/2,1}}{\pi}$  by the separation condition on  $Y$  and conclude that

$$\begin{aligned} \left( \sum_{t' \in Y} (\varphi * \varphi)((1 + \tau)^{-1/2}(t - t')) \right)^2 &= \left( \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} \left( \sum_{s \in Y_t} \varphi(x - s) \right) \varphi(x) \, dx \right)^2 \\ &< (\varphi * \varphi)(0) \cdot \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} \left| \sum_{s \in Y_t} \varphi(x - s) \right|^2 \, dx \\ &= (\varphi * \varphi)(0) \cdot \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} \sum_{s \in Y_t} |\varphi(x - s)|^2 \, dx \\ &= (\varphi * \varphi)(0) \cdot \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0) \cap \left( \bigcup_{s \in Y_t} B_{\frac{j_{d/2,1}}{2\pi}}(s) \right)} |\varphi(x)|^2 \, dx \\ &\leq (\varphi * \varphi)(0) \cdot \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} |\varphi(x)|^2 \, dx = (\varphi * \varphi)(0)^2 \end{aligned}$$

by means of Hölder's inequality and the disjointedness of  $\text{supp } \varphi(\cdot - s) = B_{\frac{j_{d/2,1}}{2\pi}}(s)$ ,  $s \in Y_t$ .

Analogously, one derives

$$\begin{aligned} \sum_{t' \in Y} \left[ \mathbb{1}_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} * \varphi \right] ((1 + \tau)^{-1/2}(t - t')) &= \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0) \cap \left( \bigcup_{s \in Y_t} B_{\frac{j_{d/2,1}}{2\pi}}(s) \right)} \varphi(x) \, dx \\ &\leq \int_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} \varphi(x) \, dx = \left[ \mathbb{1}_{B_{\frac{j_{d/2,1}}{2\pi}}(0)} * \varphi \right] (0). \end{aligned}$$

This yields  $\sum_{t' \in Y} \psi_\tau(t - t') < \psi_\tau(0)$  for  $\text{dist}(t, Y) > 0$ . As already explained in footnote 30, we have  $\psi_\tau \in C^2(\mathbb{R}^d)$ . Since  $\text{dist}(t, Y)$  is realised for at least one  $t' \in Y$ , we can estimate the sum over evaluations of  $\psi_\tau$  by bounding the remaining  $|Y| - 1$  terms by  $\psi_\tau$  evaluated at the smallest admissible radius.<sup>31</sup> This means

$$\sum_{t' \in Y} \psi_\tau(t - t') \leq \psi_\tau(\text{dist}(t, Y)e_1) + (|Y| - 1)\psi_\tau([q_\tau - \text{dist}(t, Y)]e_1) =: F(\text{dist}(t, Y))$$

<sup>31</sup>If  $\|t - t'\|_2 = \text{dist}(t, Y) \in (0, q_\tau)$  for some  $t' \in Y$ , we have  $\|t - s\|_2 \geq \|s - t'\| - \|t - t'\| \geq q_\tau - \text{dist}(t, Y)$  for any  $s \in Y$  with  $s \neq t'$ .

## 2 Condition of sparse super resolution

where  $F$  satisfies  $F(0) = \psi_\tau(0)$ ,  $F'(0) = 0$  and  $F''(0) < 0$  since  $\psi_\tau \in C^2(\mathbb{R}^d)$ . Now, we proceed similar as for (2.9) and obtain the existence of  $c_0, r_0 > 0$  such that

$$\begin{aligned} \psi_\tau(0) - \sum_{t' \in Y} \psi_\tau(t - t') &\geq F(0) - F(\text{dist}(t, Y)) \geq \frac{F''(0)}{4} \text{dist}(t, Y)^2 \\ &= \frac{4\pi^2\tau(1+\tau)^{-d/2-1}c_0}{4} \text{dist}(t, Y)^2 \end{aligned}$$

for all  $t$  with  $\text{dist}(t, Y) \leq r_0$ . Setting

$$c_1 = \max_{t \in \mathbb{R}^d: \text{dist}(t, Y) \geq r_0} \sum_{t' \in Y} (\varphi * \varphi) \left( (1+\tau)^{-1/2}(t - t') \right)$$

implies  $c_1 < (\varphi * \varphi)(0)$ . With this at hand, we bound

$$\begin{aligned} &\psi_\tau(0) - \sum_{t' \in Y} \psi_\tau(t - t') \\ &\geq 4\pi^2\tau(1+\tau)^{-d/2} \min \left[ \frac{c_0 \text{dist}(t, Y)^2}{4(1+\tau)}, (\varphi * \varphi)(0) - \sum_{t' \in Y} (\varphi * \varphi) \left( (1+\tau)^{-1/2}(t - t') \right) \right] \\ &\geq 4\pi^2\tau \left( \frac{1}{\sqrt{1+\tau}} \right)^{d+2} \min \left( \frac{c_0}{4}, \frac{(\varphi * \varphi)(0) - c_1}{\left( \frac{j_{d/2,1}}{2\pi} \right)^2} \right) \text{dist}(t, Y)^2 \end{aligned}$$

and this was the proposed result.<sup>32</sup> □

**Example 2.2.5** ( $d = 1$ ). In the univariate case, the definition (2.8) can be made more explicit. Because of  $J_{-1/2}(2\pi x) = \left( \frac{1}{\pi^2 x} \right)^{1/2} \cos(2\pi x)$  and  $j_{1/2,1} = \pi$ , see Lemma 1.3.2, one can compute

$$\varphi(x) = 1 + \sqrt{2x} \left( \frac{1}{\pi^2 x} \right)^{1/2} \cos(2\pi x) \left( \frac{2}{\pi^2} \right)^{-1/2} = 1 + \cos(2\pi x) = 2 \cos^2(\pi x).$$

This motivates the following tensor approach for an admissible function which we display in Figure 2.4. By the product structure, we have a function  $\psi$  with support in a box. This approach has the advantage to give explicit estimates on  $\psi(0) - \psi(x)$  without involving special functions.

**Lemma 2.2.6** (Support on a cube). *Let  $d \geq 1$ ,  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ , and  $\psi: \mathbb{R}^d \rightarrow \mathbb{R}$  with*

$$\varphi(x) = \begin{cases} \cos^2\left(\frac{\pi x}{q}\right), & |x| < \frac{q}{2}, \\ 0, & \text{otherwise,} \end{cases} \quad \psi = (4\pi^2 + \Delta) \bigotimes_{\ell=1}^d (\varphi * \varphi).$$

*Then,  $\psi$  is admissible if  $q \geq \sqrt{d}$  and  $d \geq 2$ . For  $q = \sqrt{d} \geq \|x\|_2$ , we have*

$$\psi(0) - \psi(x) \geq \left( \frac{3\sqrt{d}}{8} \right)^{d-1} \frac{\pi^2 \|x\|_2^2}{d^{3/2}}.$$

<sup>32</sup>We stress that we do not state that  $c_0, c_1$  in the proofs of Lemma 2.2.2 and Lemma 2.2.4 are equal. Therefore, we use  $c_d$  and  $c'_d$  respectively for the dimension dependent constants in (2.9) and (2.14).

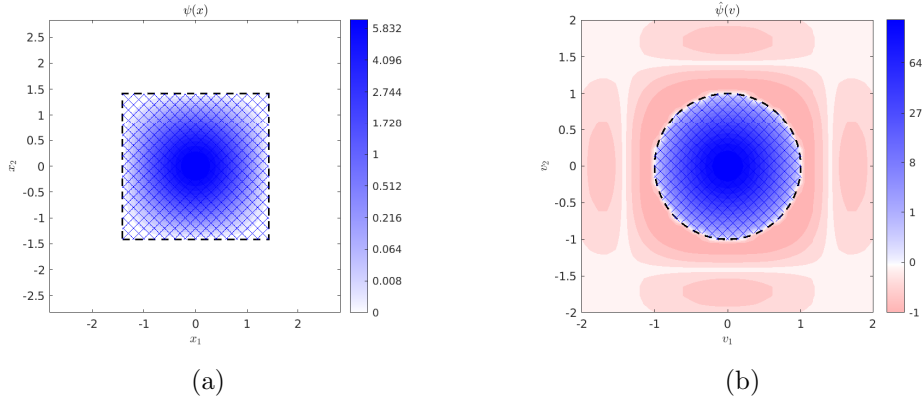


Figure 2.4: Minorant function  $\psi$  and its Fourier transform  $\hat{\psi}$  for  $\varphi$  as chosen in Lemma 2.2.6 and  $q = \sqrt{d} = \sqrt{2}$ . The function  $\psi$  is nonnegative inside of the dashed, hatched rectangle and zero outside (a) while  $\hat{\psi}$  is nonnegative inside the dash-dotted, hatched circle and nonpositive outside (b). Both functions have maximal value at zero.

*Proof.* We have  $\text{supp } \varphi = [-q/2, q/2]$ . Thus, tensorising and applying the differential operator  $\Delta$  yields  $\text{supp } \psi = [-q, q]^d$ . Moreover, we have  $\widehat{\varphi * \varphi} = (\widehat{\varphi})^2 \geq 0$  and thus  $\hat{\psi}(v) = 4\pi(1 - \|v\|_2^2) \otimes_{\ell=1}^d \widehat{\varphi * \varphi}(v)$  has the desired sign. Direct calculation gives

$$\varphi * \varphi(x) = \begin{cases} \frac{q-|x|}{4} \left(1 + \frac{1}{2} \cos\left(\frac{2\pi x}{q}\right)\right) + \frac{3}{8} \frac{q}{2\pi} \sin\left(\frac{2\pi|x|}{q}\right), & |x| < q, \\ 0, & \text{otherwise,} \end{cases}$$

and  $(\varphi * \varphi)''(x) = -\frac{4\pi^2}{q^2} \varphi * \varphi(x) + \frac{\pi^2}{q^2} \left(\frac{q}{2\pi} \sin\left(\frac{2\pi|x|}{q}\right) + q - |x|\right)$  for  $|x| < q$ . This yields

$$\psi(x) = \sum_{s=1}^d \left[ 4\pi^2 \left[1 - \frac{d}{q^2}\right] (\varphi * \varphi)(x_s) + \frac{\pi^2}{q^2} \left(\frac{\sin\left(\frac{2\pi|x_s|}{q}\right)}{2\pi/q} + q - |x_s|\right) \right] \prod_{i \neq s} (\varphi * \varphi)(x_i).$$

The global maximality of  $\psi$  in 0 for  $q \geq \sqrt{d}$  is proven by the inequality  $\sin(x) \leq x$  and the maximality of  $\varphi * \varphi$  in zero. More specifically, it is straightforward to show

$$\begin{aligned} (\varphi * \varphi)(x) &\leq \frac{3q}{8} \left(1 - \frac{|x|^2}{q^2}\right), \quad |x| \leq q, \text{ and} \\ \psi(0) - \psi(x) &\geq \psi(0) - \frac{\pi^2}{q} \sum_{s=1}^d \left[\frac{3q}{8} \left(1 - \frac{|x_s|^2}{q^2}\right)\right]^{d-1} \geq \left(\frac{3q}{8}\right)^{d-1} \frac{\pi^2 \|x\|_2^2}{q^3} \end{aligned}$$

and this was the proposed estimate.  $\square$

**Remark 2.2.7** (Cubic lower bound at critical radius for  $d = 1$ ). For  $d = 1$  and  $q = 1$  the function  $\psi$  from Lemma 2.2.6 (and from Lemma 2.2.2 with a different multiplicative constant) reads as

$$\psi(x) = \pi^2 \left(\frac{\sin(2\pi|x|)}{2\pi} + 1 - |x|\right)$$

## 2 Condition of sparse super resolution

for  $|x| \leq 1$ . The latter agrees up to the constant  $\pi^2$  with a minorant by Diederichs [37, Lem. 2.2.1], see Section 1.3. For this minorant, he realised that it does just admit a sub-optimal estimate of the form  $\psi(0) - \psi(x) \geq Cx^3$  for some  $C > 0$ . We can now understand this better by Remark 2.2.3 where we observed that  $q = 1$  is the critical radius for  $d = 1$  that does not allow the existence of an admissible function.

### 2.2.2 Ingham-type inequality for parameter difference

The following main result shows that two well-separated measures  $\mu_1 = \sum_{t \in Y^{\mu_1}} \alpha_t^{(1)} \delta_t$  and  $\mu_2 = \sum_{t \in Y^{\mu_2}} \alpha_t^{(2)} \delta_t$  with similar moments have also similar weights and nodes. We call it a *Lipschitz* result since it bounds the difference in parameters by the difference of the moments and we add the term *local* since we think at first about measures having similar moments. While main parts of the proof remain valid without the assumption of a small distance of the moment sequences and therefore allow for a global result in Subsection 2.2.3, the locality condition (2.15) ensures that the two measures have the same number of parameters. Inequalities of this kind are sometimes called *Ingham inequalities* (cf. [76, 83, 86]). The proof of our version works similar as in [37, 38] but with a completely different minorant function allowing for generalisation to  $d > 2$ , radial sampling sets and a smaller separation condition.

**Theorem 2.2.8** (Local Lipschitz). *Let  $d \geq 1$  and fix a bandlimit  $n > 0$ . We choose  $\tau > 0$  such that  $q = \frac{\sqrt{1+\tau} j_{d/2,1}}{\pi n} = \min\{\text{sep } Y^{\mu_1}, \text{sep } Y^{\mu_2}\}$  and  $\hat{\mu}_1, \hat{\mu}_2 \in \widehat{\mathcal{M}}^n(q)$ . Let  $\alpha_{\min}$  be the minimal absolute value of any weight of  $\mu_1$  or  $\mu_2$ . There exists a constant  $c_{d,\tau}^{(1)} > 0$  depending on  $d$  and  $\tau$  such that for all  $\hat{\mu}_1, \hat{\mu}_2$  with*

$$\|\hat{\mu}_1 - \hat{\mu}_2\|_2^2 < c_{d,\tau}^{(1)} \cdot n^d \alpha_{\min}^2 \quad (2.15)$$

*we find for every  $t \in Y^{\mu_1}$  exactly one  $t' = \eta(t) \in Y^{\mu_2}$  with  $\|t - \eta(t)\|_{\mathbb{T}^d} < \frac{q}{2}$  and vice versa. Moreover, there are  $c_{d,\tau}^{(2)}, c_{d,\tau}^{(3)} > 0$  such that (2.15) also implies*

$$\|\hat{\mu}_1 - \hat{\mu}_2\|_2^2 \geq \sum_{t \in Y^{\mu_1}} c_{d,\tau}^{(2)} n^{d+2} \alpha_{\min}^2 \|t - \eta(t)\|_{\mathbb{T}^d}^2 + c_{d,\tau}^{(3)} n^d |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2. \quad (2.16)$$

*For  $d \geq 2$  and a larger separation  $\min(\text{sep } Y^{\mu_1}, \text{sep } Y^{\mu_2}) \geq \frac{2d}{n}$  one can take*

$$c_{d,\tau}^{(1)} = \left(\frac{3}{2}\right)^{d-1} d^{-d/2}, \quad c_{d,\tau}^{(2)} = \frac{c_{d,\tau}^{(1)}}{2d^2}, \quad \text{and} \quad c_{d,\tau}^{(3)} = \frac{c_{d,\tau}^{(1)}}{4}.$$

*Proof.* Due to the condition on the separation of  $\hat{\mu}_1$  and  $\hat{\mu}_2$ , we know that for every  $t \in Y^{\mu_1}$  there is at most one  $\eta(t) \in Y^{\mu_2}$  with  $\|t - \eta(t)\|_{\mathbb{T}^d} < \frac{q}{2}$ .<sup>33</sup> As visualised in Figure 2.5, we decompose the joint node set  $Y := Y^{\mu_1} \cup Y^{\mu_2}$  into  $Y_1 \subset Y^{\mu_1}$ ,  $Y_2 \subset Y^{\mu_2}$  and  $Y_3 \subset Y^{\mu_1} \cup Y^{\mu_2}$ , see also [38, Thm. 3.6], with:

- (i)  $Y_3 := \{t \in Y : \text{For all } t' \in Y \text{ with } t \neq t' \text{ one has } \|t - t'\|_{\mathbb{T}^d} \geq \frac{q}{2}\}$
- (ii) For all  $t \in Y_1$  there is exactly one  $\eta(t) \in Y_2$  with  $\|t - \eta(t)\|_{\mathbb{T}^d} < \frac{q}{2}$ .

<sup>33</sup>If there were  $t \in Y^{\mu_1}$ ,  $t_1, t_2 \in Y^{\mu_2}$ ,  $t_1 \neq t_2$  with  $\|t - t_1\|_{\mathbb{T}^d} < \frac{q}{2}$  and  $\|t - t_2\|_{\mathbb{T}^d} < \frac{q}{2}$ , we have  $\|t_1 - t_2\|_{\mathbb{T}^d} < q$  which is a contradiction to the assumption  $\text{sep } Y^{\mu_2} \geq q$ .

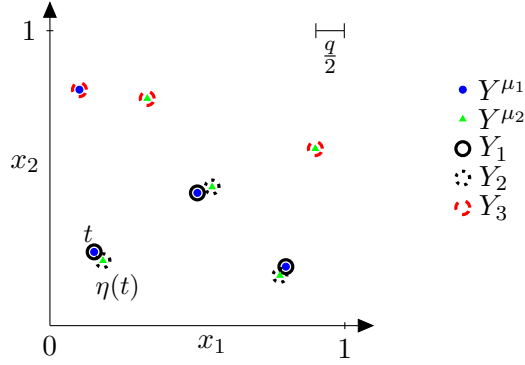


Figure 2.5: Visualisation of the sets  $Y^{\mu_1}$  (blue circles) and  $Y^{\mu_2}$  (green triangles) as well as their subsets  $Y_1 \subset Y^{\mu_1}$  (densely black circles) and  $Y_2 \subset Y^{\mu_2}$  (dotted black circles). Nodes without a neighbour closer than  $\frac{q}{2}$  belong to  $Y_3$  (dashed red circles).

The function  $\psi_{\tau,n}(x) := n^d \psi_{\tau}(n \cdot x)$  with  $\psi_{\tau}$  from Lemma 2.2.2 has compact support in  $B_q(0)$  and its Fourier transform  $\hat{\psi}_{\tau,n}$  decreases fast enough in order to apply the Poisson summation formula.<sup>34</sup> Together with the conventions

$$\tilde{\alpha}_t = \begin{cases} \alpha_t^{(1)}, & t \in Y^{\mu_1}, \\ -\alpha_t^{(2)}, & t \in Y^{\mu_2}, \end{cases} \quad \text{and} \quad \tilde{Y}_t := \begin{cases} Y^{\mu_2}, & t \in Y^{\mu_1}, \\ Y^{\mu_1}, & t \in Y^{\mu_2}, \end{cases} \quad (2.17)$$

this gives

$$\begin{aligned} \hat{\psi}_{\tau,n}(0) \sum_{\substack{k \in \mathbb{Z}^d \\ \|k\|_2 \leq n}} |\hat{\mu}_1(k) - \hat{\mu}_2(k)|^2 &\geq \sum_{k \in \mathbb{Z}^d} |\hat{\mu}_1(k) - \hat{\mu}_2(k)|^2 \hat{\psi}_{\tau,n}(k) \\ &= \sum_{t,t' \in Y} \tilde{\alpha}_t \overline{\tilde{\alpha}_{t'}} \sum_{\ell \in \mathbb{Z}^d} \psi_{\tau,n}(t - t' + \ell). \end{aligned} \quad (2.18)$$

For every  $t, t' \in Y$  there is a unique  $\ell_{t,t'} \in \mathbb{Z}^d$  such that  $t - t' + \ell_{t,t'} \in [-\frac{1}{2}, \frac{1}{2}]^d$  and because of  $\text{supp } \psi_{\tau,n} = B_q(0)$  we have  $\sum_{\ell \in \mathbb{Z}^d} \psi_{\tau,n}(t - t' + \ell) = \psi_{\tau,n}(t - t' + \ell_{t,t'})$ . Defining for each  $t \in Y^{\mu_1}$  the shifted node set by  $Y_t^{\hat{\mu}_2} := \{t' - \ell_{t,t'} : t' \in Y^{\mu_2}\}$  preserves the separation in the wrap around distance, as  $Y_t^{\hat{\mu}_2} = Y^{\mu_2}$  on  $\mathbb{T}^d$ . For  $t' \in Y^{\mu_2}$  we analogously define  $Y_t^{\hat{\mu}_1}$ . By the  $q$ -separation of  $\mu_1, \mu_2$ , equation (2.18) rewrites as

$$\begin{aligned} &\psi_{\tau,n}(0) \left( \sum_{t \in Y^{\mu_1}} |\alpha_t^{(1)}|^2 + \sum_{t' \in Y^{\mu_2}} |\alpha_{t'}^{(2)}|^2 \right) - \sum_{t \in Y^{\mu_1}} \sum_{t' \in Y_t^{\hat{\mu}_2}} \psi_{\tau,n}(t - t') 2\Re \left( \alpha_t^{(1)} \overline{\alpha_{t'}^{(2)}} \right) \\ &= \sum_{t \in Y^{\mu_1}} \sum_{t' \in Y_t^{\hat{\mu}_2}} \psi_{\tau,n}(t - t') |\alpha_t^{(1)} - \alpha_{t'}^{(2)}|^2 + \sum_{t \in Y^{\mu_1}} |\alpha_t^{(1)}|^2 \left[ \psi_{\tau,n}(0) - \sum_{t' \in Y_t^{\hat{\mu}_2}} \psi_{\tau,n}(t - t') \right] \\ &\quad + \sum_{t' \in Y^{\mu_2}} |\alpha_{t'}^{(2)}|^2 \left[ \psi_{\tau,n}(0) - \sum_{t \in Y_t^{\hat{\mu}_1}} \psi_{\tau,n}(t - t') \right] \end{aligned}$$

<sup>34</sup>For the Poisson summation formula from Theorem 1.2.3 it is sufficient that both the function and its Fourier transform decay with the rate  $\mathcal{O}(\|x\|_2^{-d-\epsilon})$  and  $\mathcal{O}(\|v\|_2^{-d-\epsilon})$  respectively. The decay of  $\psi_{\tau,n}$  is trivially fast enough by the compact support and for the decay of its Fourier transform see footnote 30.

## 2 Condition of sparse super resolution

$$\begin{aligned}
&\geq \sum_{t \in Y_1} \psi_{\tau,n}(t - \eta(t)) |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2 + \sum_{t \in Y_3} |\tilde{\alpha}_t|^2 \left[ \psi_{\tau,n}(0) - \sum_{t' \in \tilde{Y}_t} \psi_{\tau,n}(t - t') \right] \\
&\quad + \sum_{t \in Y_1} \left( |\alpha_t^{(1)}|^2 + |\alpha_{\eta(t)}^{(2)}|^2 \right) c'_d \tau (1 + \tau)^{-d/2-1} n^{d+2} \|t - \eta(t)\|_{\mathbb{T}^d}^2
\end{aligned} \tag{2.19}$$

where the last inequality is due to (2.14). For  $t \in Y_3$  we have by definition of  $Y_3$  that  $\text{dist}(t, \tilde{Y}_t) \geq \frac{q}{2}$ . By Lemma 2.2.4, we therefore obtain

$$\begin{aligned}
\sum_{\substack{k \in \mathbb{Z}^d \\ \|k\|_2 \leq n}} |\hat{\mu}_1(k) - \hat{\mu}_2(k)|^2 &\geq \sum_{t \in Y_3} \frac{\psi_{\tau,n}(0) - \sum_{t' \in \tilde{Y}_t} \psi_{\tau,n}(t - t')}{\hat{\psi}_{\tau,n}(0)} |\tilde{\alpha}_t|^2 \\
&\geq \sum_{t \in Y_3} \frac{\min \left( c'_d \tau (1 + \tau)^{-d/2-1} \left[ n \frac{q}{2} \right]^2, \psi_{\tau}(0) \right)}{\hat{\psi}_{\tau}(0)} n^d \alpha_{\min}^2 \\
&\geq \frac{\min \left( c'_d \tau (1 + \tau)^{-d/2-1} \frac{j_{d/2,1}^2}{4\pi^2}, \psi_{\tau}(0) \right)}{\hat{\psi}_{\tau}(0)} n^d \alpha_{\min}^2 |Y_3|.
\end{aligned} \tag{2.20}$$

If now (2.15) holds with  $c_{d,\tau}^{(1)} := \min \left( c'_d \tau (1 + \tau)^{-d/2-1} \frac{j_{d/2,1}^2}{4\pi^2}, \psi_{\tau}(0) \right) / \hat{\psi}_{\tau}(0)$  this yields  $|Y_3| < 1$  meaning  $Y_3 = \emptyset$ . So we know already that for all  $t \in Y_1 = Y^{\mu_1}$  there is  $t' = \eta(t) \in Y_2 = Y^{\mu_2}$  with  $\|t - \eta(t)\|_{\mathbb{T}^d} < \frac{q}{2}$ . Our previous estimates can then be summarised to (2.16) with  $c_{d,\tau}^{(2)} := 2c'_d \tau (1 + \tau)^{-d/2-1} / \hat{\psi}_{\tau}(0)$  and  $c_{d,\tau}^{(3)} := \psi_{\tau}(n \frac{q}{2} e_1) / \hat{\psi}_{\tau}(0)$ . For the last part of the statement, one can redo the proof with the function  $n^d \psi(n \cdot x)$  where  $\psi$  as in Lemma 2.2.6 in order to obtain these constants explicitly under a stronger separation condition.  $\square$

**Remark 2.2.9** (Optimal orders). In the univariate case, Diederichs [37, Lemma 2.24] gives simple examples that the orders in  $n, \alpha_{\min}, \|t - \eta(t)\|_{\mathbb{T}^d}$ , and  $|\alpha_t - \alpha_{\eta(t)}|$  are optimal in (2.15) and (2.16). Since  $|Y^{\mu_1}| = |Y^{\mu_2}| = 1$  in these examples, they directly carry over to the  $d$ -variate case. In fact, one can choose any  $t, t' \in \mathbb{T}^d$  and  $\alpha_{\min} > 0$  to calculate

$$\begin{aligned}
\|\widehat{\alpha_{\min} \delta_t} - \widehat{\alpha_{\min} \delta_{t'}}\|_2^2 &= \sum_{k: \|k\|_2 \leq n} \alpha_{\min}^2 |1 - e^{2\pi i(t-t')k}|^2 \\
&= 4\alpha_{\min}^2 \sum_{k: \|k\|_2 \leq n} |\sin(\pi k(t - t'))|^2 \\
&\leq 4\pi^2 \alpha_{\min}^2 \|t - t'\|_{\mathbb{T}^d}^2 \sum_{k: \|k\|_2 \leq n} \|k\|_2^2 \\
&\leq 4\pi^2 \alpha_{\min}^2 \|t - t'\|_{\mathbb{T}^d}^2 2^d n^{d+2}
\end{aligned}$$

as well as

$$\|\widehat{\alpha_1 \delta_t} - \widehat{\alpha_2 \delta_t}\|_2^2 = \sum_{k: \|k\|_2 \leq n} |\alpha_1 - \alpha_2|^2 \leq 2^d n^d |\alpha_1 - \alpha_2|^2$$

for any  $\alpha_1, \alpha_2 \in \mathbb{C}$ . This shows that the dependency in  $n, \alpha_{\min}, \|t - \eta(t)\|$ , and  $|\alpha_t - \alpha_{\eta(t)}|$  as presented in (2.16) is optimal in general. Moreover, taking  $t' \in \mathbb{T}^d$  such that  $\mu_2 =$

$\mu_1 + \alpha_{\min} \delta_{t'}$  satisfies  $\mu_1, \mu_2 \in \mathcal{M}(q)$  gives two measures with unequal cardinality of the node set but with moments satisfying

$$\|\hat{\mu}_1 - \hat{\mu}_2\|_2^2 = \alpha_{\min}^2 \sum_{k: \|k\|_2 \leq n} 1 \leq \frac{2\pi^{d/2}}{\Gamma(d/2)} \alpha_{\min}^2 \int_0^{n+1} r^{d-1} dr = \frac{2\pi^{d/2}}{d \cdot \Gamma(d/2)} (n+1)^d \alpha_{\min}^2.$$

This yields optimality of the orders of  $\alpha_{\min}$  and  $n$  in (2.15). We do not propose optimality of the constants  $c_{d,\tau}^{(1)}$ ,  $c_{d,\tau}^{(2)}$  and  $c_{d,\tau}^{(3)}$ . However, we are interested in their dependency on  $\tau > 0$ . Therefore, we compute  $\hat{\psi}_\tau(0) = 4\pi^2(1+\tau)\hat{\varphi}(0)^2$  and

$$\begin{aligned} \psi_\tau(0) &= 4\pi^2(1+\tau) \int_{\mathbb{R}^d} (1 - \|v\|_2^2) [\hat{\varphi}(\sqrt{1+\tau}v)]^2 dv \\ &\geq 4\pi^2 \int_{\mathbb{R}^d} (1 - \|\sqrt{1+\tau}v\|_2^2) [\hat{\varphi}(\sqrt{1+\tau}v)]^2 dv \\ &= (1+\tau)^{-d/2} 4\pi^2 \int_{\mathbb{R}^d} (1 - \|w\|_2^2) [\hat{\varphi}(w)]^2 dw. \end{aligned}$$

Additionally, one finds

$$\psi_\tau\left(\frac{nq}{2}e_1\right) = (1+\tau)^{-d/2} \left[ \psi_0\left(\frac{1}{2}e_1\right) + \tau(\varphi * \varphi)\left(\frac{1}{2}e_1\right) \right] \geq (1+\tau)^{-d/2} \psi_0\left(\frac{1}{2}e_1\right)$$

and this yields

$$\begin{aligned} c_{d,\tau}^{(1)} &= \min\left(c'_d \tau (1+\tau)^{-d/2-1} \frac{j_{d/2,1}^2}{4\pi^2}, \psi_\tau(0)\right) / \hat{\psi}_\tau(0) \gtrsim_d \tau (1+\tau)^{-d/2-2}, \\ c_{d,\tau}^{(2)} &= 2c'_d \tau (1+\tau)^{-d/2-1} / \hat{\psi}_\tau(0) \gtrsim_d \tau (1+\tau)^{-d/2-2} \text{ and} \\ c_{d,\tau}^{(3)} &= \psi_\tau\left(\frac{nq}{2}e_1\right) / \hat{\psi}_\tau(0) \gtrsim_d (1+\tau)^{-d/2-2}. \end{aligned}$$

This analysis shows that none of the three constants goes to zero faster than linear in  $\tau$  if one approaches the critical separation, i.e. in the limit  $\tau \rightarrow 0$ .

**Remark 2.2.10** (Critical separation). For a statement like Theorem 2.2.8, one cannot hope to reduce the need of separation below  $\frac{j_{d/2,1}}{\pi n}$  using the approach with a minorising function as there is no admissible function in this case, see Remark 2.2.3. Note that the univariate case  $d = 1$  is not included in the second part of the theorem because  $\frac{j_{1/2,1}}{\pi n} = \frac{1}{n} = \frac{d}{n}$  is equal and not larger than the critical radius. This explains why we and also others before (cf. [37, 38]) have seen a cubic rate being worse than the optimal quadratic rate in  $\|t - \eta(t)\|_{\mathbb{T}^d}$ .<sup>35</sup> Moreover, we emphasise that both separation conditions in Theorem 2.2.8 have the same order in  $d$  because we can compute  $\frac{j_{d/2,1}}{\pi} = \frac{d}{2\pi} + \frac{1.855757}{2^{-1/3}\pi} d^{1/3} + \mathcal{O}(d^{-1/3})$  as  $d \rightarrow \infty$  by Lemma 1.3.2. On the other hand, this linear dependency in  $d$  cannot be considerably improved. This can be seen by analysing the simpler problem of a uniformly bounded smallest singular value of the corresponding Vandermonde matrix. It was observed in [88, last paragraph] that this smallest singular value scaled by  $n^{-d/2}$  goes to zero as  $d \rightarrow \infty$  if the separation  $q$  of the set satisfies  $nq = o(d)$ . Consequently, a linear rate in  $d$  for the separation meets our expectations.

<sup>35</sup>We have observed the cubic rate already in Remark 2.2.7.

## 2 Condition of sparse super resolution

**Remark 2.2.11** (Discrete and continuous data). Many publications including for example [12] assume access to continuous Fourier data, i.e.  $\hat{\mu}(k)$  for  $k \in B_n(0)$ . In contrast to this, we think that in practical implementations using the FFT one obtains only discrete moments  $\hat{\mu}(k)$  where  $k \in B_n(0) \cap \mathbb{Z}^d$ . Hence, our considerations in Theorem 2.2.8 dealt with vectors of Fourier moments. Despite of that, it should be mentioned at this point that all computations in the proof of Theorem 2.2.8 can also be made by replacing the discrete sums by integrals and then applying the inverse Fourier transform instead of the Poisson summation formula. Therefore, assuming discrete data is no restriction and all results following in this chapter can also be made for continuous inputs.

### 2.2.3 Condition number and diffraction limit

As an immediate consequence of Theorem 2.2.8 we can deduce that  $\mu_1, \mu_2 \in \mathcal{M}(q)$  with equal moments in  $\widehat{\mathcal{M}}^n(q)$ ,  $n \cdot q > \frac{j_{d/2,1}}{\pi}$ , must be equal since the differences in the parameters are bounded by the difference of the Fourier coefficients, see (2.16). Therefore, each moment sequence  $\hat{\mu} \in \widehat{\mathcal{M}}^n(q)$  with  $n \cdot q > \frac{j_{d/2,1}}{\pi}$  can be mapped uniquely to  $\mu \in \mathcal{M}(q)$ .

**Definition 2.2.12** (Multivariate reconstruction map). Let  $n \in \mathbb{N}$ ,  $d \geq 1$ . Assume we are interested in the reconstruction of a measure  $\mu$  with  $q$ -separated nodes given its moments  $\hat{\mu}(k)$  for  $k$  in the sampling set  $\{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}$  and  $n \cdot q > \frac{j_{d/2,1}}{\pi}$ . Then, the *multivariate reconstruction map* is

$$\mathcal{R} : \widehat{\mathcal{M}}^n(q) \rightarrow \mathcal{M}(q), \quad \left( \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k} \right)_{k \in \mathbb{Z}^d : \|k\|_2 \leq n} \mapsto \sum_{t \in Y} \alpha_t \delta_t.$$

We want to study the condition of this mapping in order to quantify the condition of super resolution and to understand the diffraction limit. This means that we need to bound differences in the space of Fourier moments by differences in the parameter space  $\mathcal{M}(q)$  and without the local condition (2.15) we can view such a result as a *global Lipschitz* result for the recovery map  $\mathcal{R}$ . We fix the  $\ell^2$ -norm on the space of moments  $\widehat{\mathcal{M}}^n(q)$  and choose the Wasserstein norm from Section 1.4 for  $\mathcal{M}(q)$ . Considering this metric for the parameter space has two main advantages. First, it links information about nodes and weights in one term. Secondly, the Wasserstein distance allows to compare parameter sets whose cardinality is not necessarily equal. Beyond that, we emphasise that Theorem 2.2.13 is completely independent of the minimal absolute value of the weights.

**Theorem 2.2.13** (Global Lipschitz). *There exist constants  $c_{d,\tau}^{(4)}, c_{d,\tau}^{(5)} > 0$  such that the reconstruction map  $\mathcal{R}$  satisfies*

$$W_1(\mathcal{R}(\hat{\mu}_1), \mathcal{R}(\hat{\mu}_2)) \leq \frac{c_{d,\tau}^{(4)} \sqrt{M}}{n^{d/2}} \|\hat{\mu}_1 - \hat{\mu}_2\|_2 \leq c_{d,\tau}^{(5)} \|\hat{\mu}_1 - \hat{\mu}_2\|_2$$

for all  $\hat{\mu}_1, \hat{\mu}_2 \in \widehat{\mathcal{M}}^n \left( \sqrt{1 + \tau} \frac{j_{d/2,1}}{\pi n} \right)$  if we restrict the reconstruction map to node sets with cardinality at most  $M$ . Under the stronger separation condition  $\hat{\mu}_1, \hat{\mu}_2 \in \widehat{\mathcal{M}}^n \left( \frac{2d}{n} \right)$ ,  $d \geq 2$ , we can explicitly derive  $c_{d,\tau}^{(4)} = \frac{1}{2} \sqrt{\frac{3d}{d^{d/2}}} \left( \frac{2}{3} \right)^{(d-1)/2}$  and  $c_{d,\tau}^{(5)} = \frac{1}{2}$ , i.e.

$$W_1(\mathcal{R}(\hat{\mu}_1), \mathcal{R}(\hat{\mu}_2)) \leq \frac{1}{2} \|\hat{\mu}_1 - \hat{\mu}_2\|_2.$$



*Proof.* For  $t \in Y_1$  we define  $\boldsymbol{\alpha}_t = (\alpha_t^{(1)}, \alpha_{\eta(t)}^{(2)})^\top \in \mathbb{C}^2$  and bound

$$\begin{aligned}
 \frac{4}{\sqrt{dM}} W_1(\mu_1, \mu_2) &= \frac{4}{\sqrt{dM}} \sup_{\substack{f: \text{Lip}(f) \leq 1, \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}}} \left| \int_{\mathbb{T}^d} f(x) d(\mu_1 - \mu_2)(x) \right| \\
 &= \frac{4}{\sqrt{dM}} \sup_{\substack{f: \text{Lip}(f) \leq 1, \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}}} \left| \sum_{t \in Y_3} \tilde{\alpha}_t f(t) + \sum_{t \in Y_1} \alpha_t^{(1)} f(t) - \alpha_{\eta(t)}^{(2)} f(\eta(t)) \right| \\
 &\leq \frac{4}{\sqrt{dM}} \sup_{\substack{f: \text{Lip}(f) \leq 1, \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}}} \left| \sum_{t \in Y_3} \tilde{\alpha}_t f(t) \right| + \left| \sum_{t \in Y_1} \alpha_t^{(1)} (f(t) - f(\eta(t))) + (\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}) f(\eta(t)) \right| \\
 &\leq \left[ \sum_{t \in Y_3} |\tilde{\alpha}_t|^2 \right]^{1/2} + \frac{4}{\sqrt{d}} \left[ \sum_{t \in Y_1} \|\boldsymbol{\alpha}_t\|_2^2 \|t - \eta(t)\|_{\mathbb{T}^d}^2 \right]^{1/2} + \left[ \sum_{t \in Y_1} |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2 \right]^{1/2}
 \end{aligned}$$

using the notation of the proof of Theorem 2.2.8. Then, the inequality  $(a + b + c)^2 \leq 3a^2 + 3b^2 + 3c^2$  for  $a, b, c \in \mathbb{R}$  allows to derive

$$\frac{16W_1(\mu_1, \mu_2)^2}{3dM} \leq \sum_{t \in Y_3} |\tilde{\alpha}_t|^2 + \frac{16}{d} \sum_{t \in Y_1} \|\boldsymbol{\alpha}_t\|_2^2 \|t - \eta(t)\|_{\mathbb{T}^d}^2 + \sum_{t \in Y_1} |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2. \quad (2.21)$$

Using (2.19) and (2.20), we bound the difference in the moments from below by

$$\begin{aligned}
 \|\hat{\mu}_1 - \hat{\mu}_2\|_2^2 &\geq c_{d,\tau}^{(1)} n^d \left[ \sum_{t \in Y_3} |\tilde{\alpha}_t|^2 + \frac{c_{d,\tau}^{(2)} n^2}{2c_{d,\tau}^{(1)}} \sum_{t \in Y_1} \|t - \eta(t)\|_{\mathbb{T}^d}^2 \|\boldsymbol{\alpha}_t\|_2^2 + \frac{c_{d,\tau}^{(3)}}{c_{d,\tau}^{(1)}} |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2 \right] \\
 &\geq c_{d,\tau}^{(1)} n^d \min \left( 1, \frac{c_{d,\tau}^{(3)}}{c_{d,\tau}^{(1)}}, \frac{dc_{d,\tau}^{(2)} n^2}{32c_{d,\tau}^{(1)}} \right) \frac{16W_1(\mu_1, \mu_2)^2}{3dM}
 \end{aligned}$$

where we applied (2.21). If  $\hat{\mu}_1, \hat{\mu}_2 \in \widehat{\mathcal{M}}^n \left( \frac{2d}{n} \right)$  we can use the constants from Theorem 2.2.8 to calculate

$$\frac{c_{d,\tau}^{(3)}}{c_{d,\tau}^{(1)}} = \frac{1}{4} \quad \text{and} \quad \frac{dc_{d,\tau}^{(2)} n^2}{32c_{d,\tau}^{(1)}} = \frac{n^2}{64d} \geq \frac{4d^2}{64d(\text{sep } Y^{\mu_1})^2} \geq \frac{d}{4} \geq \frac{1}{2}.$$

This gives

$$c_{d,\tau}^{(4)} = \sqrt{\frac{3d}{4c_{d,\tau}^{(1)}}} = \frac{1}{2} \sqrt{3dd^{d/2}} \left( \frac{2}{3} \right)^{(d-1)/2}.$$

Finally, the global Lipschitz estimate is further simplified by noting that  $M \leq \left( n/(2\sqrt{d}) \right)^d$  through the separation condition of  $Y^{\mu_1}, Y^{\mu_2} \subset \mathbb{T}^d$ . We then end up with

$$c_{d,\tau}^{(5)} = \frac{c_{d,\tau}^{(4)}}{2^{d/2} d^{d/4}} = \frac{3\sqrt{d}}{2\sqrt{2}} \left( \frac{1}{3} \right)^{d/2} \leq \frac{3\sqrt{2}}{2\sqrt{2}} \left( \frac{1}{3} \right)^{2/2} = \frac{1}{2}$$

as the latter expression in  $d \geq 2$  becomes maximal for  $d = 2$ .  $\square$

## 2 Condition of sparse super resolution

While the reconstruction map  $\mathcal{R}$  is just defined on exact data from  $\widehat{\mathcal{M}}^n \left( \sqrt{1 + \tau} \frac{j_{d/2,1}}{\pi n} \right)$ , we are more interested in perturbation results that hold also for perturbed inputs. As already observed by Diederichs in [37], a Lipschitz result like Theorem 2.2.8 or Theorem 2.2.13 allows to conclude a-posteriori error bounds. We formulate the following perturbation result for the best sparse solution if we are given noisy measurements.

**Corollary 2.2.14** (Perturbation result). *Let  $\varrho > 0, \tau > 0$ . Assume that the measure  $\mu_0 \in \mathcal{M} \left( \sqrt{1 + \tau} \frac{j_{d/2,1}}{\pi n} \right)$  is  $M$ -sparse and that one has access to its noisy Fourier coefficients (with frequencies  $k \in \mathbb{Z}^d, \|k\|_2 \leq n$ )*

$$\hat{\mu} = \widehat{\mu}_0 + \hat{\rho}, \quad \|\hat{\rho}\|_2 \leq \varrho.$$

Then, one has  $W_1(\mu_0, \nu_*) \leq 2c_{d,\tau}^{(5)} \varrho$  for any at most  $M$ -sparse best approximation

$$\nu_* \in \underset{\substack{\nu \in \mathcal{M} \left( \frac{2d}{n} \right) \\ |Y^\nu| \leq M}}{\operatorname{argmin}} \|\hat{\mu} - \hat{\nu}\|_2. \quad (2.22)$$

*Proof.* The existence of a best approximation defined by (2.22) follows from the fact that we can parameterise the set  $\{\nu \in \mathcal{M} \left( \frac{2d}{n} \right) : |Y^\nu| \leq M\}$  as the image of a closed subset of  $(\mathbb{T}^d \times \mathbb{C})^M$  under a continuous mapping. As any best approximation  $\nu_*$  which might not necessarily be unique satisfies the conditions of Theorem 2.2.13, we get

$$W_1(\mu_0, \nu_*) \leq \frac{c_{d,\tau}^{(4)} \sqrt{M}}{n^{d/2}} \|\widehat{\mu}_0 - \widehat{\nu}_*\|_2 \leq \frac{c_{d,\tau}^{(4)} \sqrt{M}}{n^{d/2}} (\varrho + \|\hat{\mu} - \widehat{\nu}_*\|_2)$$

by the triangle inequality. The last inequality is due to the fact that  $\mu_0$  is admissible to the optimisation problem (2.22).  $\square$

With this perturbation result at hand, we can now come to the definition of a condition number for the super resolution problem. While we defined a *structured* absolute condition number for  $\mathcal{R}$  at  $\hat{\mu} \in \widehat{\mathcal{M}}^n(q)$  by

$$\kappa_{\text{abs}}^{\text{str}}(\hat{\mu}, q, n, d) = \lim_{\varrho \rightarrow 0} \sup_{\substack{\|\hat{\rho}\|_2 \leq \varrho \\ \hat{\mu} + \hat{\rho} \in \widehat{\mathcal{M}}^n(q)}} \frac{W_1(\mathcal{R}(\hat{\mu} + \hat{\rho}), \mathcal{R}(\hat{\mu}))}{\|\hat{\rho}\|_2}$$

in [68], a meaningful definition from the computational point of view would incorporate unstructured perturbations  $\hat{\rho}$ .<sup>36</sup> A generalisation to an unstructured absolute or relative condition number requires an extension of the reconstruction map  $\mathcal{R}$  to a neighbourhood of  $\widehat{\mathcal{M}}^n \left( \frac{2d}{n} \right)$  which is not straightforward as it is not obvious whether (2.22) is unique. We circumvent this problem by allowing a set of possible reconstructions as outputs of the reconstruction map. Any optimal black box algorithm can then just select the best out of the set of possible reconstructions.<sup>37</sup>

<sup>36</sup>Note that the unstructured condition number satisfies due to Theorem 2.2.13  $\kappa_{\text{abs}}^{\text{str}}(\hat{\mu}, n, q, d) \leq \frac{1}{2}$  if  $q > \frac{2d}{n}$ . By the global Lipschitz result, the bound on the norm of  $\hat{\rho}$  in terms of  $\varrho$  is not necessary for the existence of the supremum and can also be dropped. A small  $\varrho$  just assures that  $\mathcal{R}(\hat{\mu} + \hat{\rho})$  and  $\mathcal{R}(\hat{\mu})$  have the same number of parameters by the locality condition in Theorem 2.2.8. Therefore, we will omit to take a limit  $\varrho \rightarrow 0$ .

<sup>37</sup>The idea to study set-valued outcomes for the analysis of the condition of approximation problems can also be found in [18, 40] and the references therein.

**Definition 2.2.15** (Set-valued reconstruction map). Let  $n \in \mathbb{N}$ ,  $d \geq 1$  and  $q > 0$ . Then, the reconstruction map  $\mathcal{R}$  can be extended to perturbed moment sequences by

$$\mathcal{R} : \mathbb{C}^{|\mathcal{I}|} \rightarrow \mathcal{P}(\mathcal{M}(q)), \quad \hat{\mu} \mapsto \operatorname{argmin}_{\nu \in \mathcal{M}(q)} \|\hat{\mu} - \nu\|_2 \quad (2.23)$$

where  $|\mathcal{I}| = |\{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}|$  and  $\mathcal{P}(\mathcal{M}(q))$  denotes the power set of  $\mathcal{M}(q)$ .

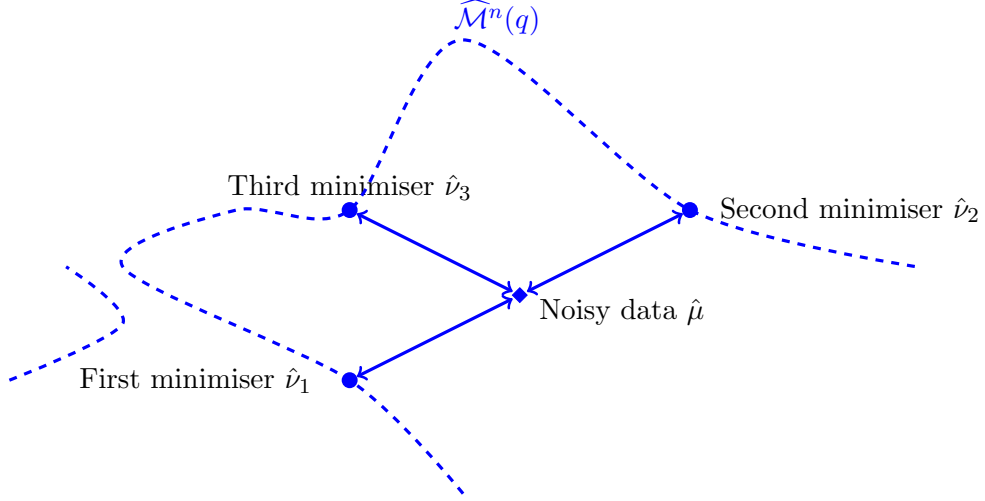


Figure 2.6: Sketch of the set-valued reconstruction map which allows to obtain multiple minimisers having equal distance to the noisy data. In this example, we would have  $\mathcal{R}(\hat{\mu}) = \{\nu_1, \nu_2, \nu_3\}$ .

The definition through the argmin is well-defined as the selection of a separation  $q$  bounds the number of possible parameters and thus  $\mathcal{M}(q)$  is closed. Moreover, this definition extends  $\mathcal{R}$  to  $q < \frac{j_{1,1}}{\pi n}$  where the mapping of measures to moments is not necessarily injective. We sketch the extension of  $\mathcal{R}$  to  $\mathbb{C}^{|\mathcal{I}|}$  in Figure 2.6. Based on this, we define the condition number of super resolution as the worst unstructured absolute condition at any  $\hat{\mu} \in \widehat{\mathcal{M}}^n(q)$ . The ratio behind this is to consider the difference between the worst possible ground truth and its best reconstruction via the least square problem in the presence of noise.<sup>38</sup>

**Definition 2.2.16** (Condition number of super resolution). Take again  $d, n, M \in \mathbb{N}$  and  $q > 0$ . Then, we define the *condition number of super resolution* as

$$\kappa_{\text{abs}}(q, n, d, M) := \sup_{\substack{\hat{\mu} \in \widehat{\mathcal{M}}^n(q) \\ |Y^{\hat{\mu}}| \leq M}} \sup_{\substack{\hat{\rho} \in \mathbb{C}^{|\mathcal{I}|} \\ \hat{\rho} \neq 0}} \inf_{\nu \in \mathcal{R}(\hat{\mu} + \hat{\rho})} \frac{W_1(\nu, \mu)}{\|\hat{\rho}\|_2}.$$

**Corollary 2.2.17** (Condition number for well-separated nodes). *If the separation fulfils  $nq = \sqrt{1 + \tau} \frac{j_{d/2,1}}{\pi}$  with  $\tau > 0$ , i.e. for well separated nodes, we have*

$$\kappa_{\text{abs}}(q, n, d, M) \leq 2c_{d,\tau}^{(5)}$$

<sup>38</sup>We do not claim that this is the only possible way to define the condition number of super resolution. For instance, Breiding and Vannieuwenhoven [18] study the condition of approximation problems like (2.23) in a different way.

## 2 Condition of sparse super resolution

with the constant  $c_{d,\tau}^{(5)}$  as specified in the proof of Theorem 2.2.13. We remark that

$$c_{d,\tau}^{(5)} \lesssim_d \frac{1+\tau}{\sqrt{\tau}}$$

such that the bound does not increase too heavily if one approaches  $\tau \rightarrow 0$ .

*Proof.* As before, we set  $q_\tau = \sqrt{1+\tau} \frac{j_{d/2,1}}{\pi}$  and apply Theorem 2.2.13 in order to bound

$$\begin{aligned} \kappa_{\text{abs}}(q, n, d, M) &\leq \sup_{\substack{\hat{\mu} \in \widehat{\mathcal{M}}^n(q) \\ |Y^\mu| \leq M}} \sup_{\substack{\hat{\rho} \in \mathbb{C}^{|\mathcal{I}|} \\ \hat{\rho} \neq 0}} \inf_{\nu \in \mathcal{R}(\hat{\mu} + \hat{\rho})} c_{d,\tau}^{(5)} \frac{\|\hat{\nu} - \hat{\mu}\|_2}{\|\hat{\rho}\|_2} \\ &\leq \sup_{\substack{\hat{\mu} \in \widehat{\mathcal{M}}^n(q) \\ |Y^\mu| \leq M}} \sup_{\substack{\hat{\rho} \in \mathbb{C}^{|\mathcal{I}|} \\ \hat{\rho} \neq 0}} \inf_{\nu \in \mathcal{R}(\hat{\mu} + \hat{\rho})} c_{d,\tau}^{(5)} \frac{\|\hat{\nu} - \hat{\mu} - \hat{\rho}\|_2 + \|\hat{\rho}\|_2}{\|\hat{\rho}\|_2} \\ &\leq 2c_{d,\tau}^{(5)} \end{aligned}$$

where the last inequality is valid since  $\mu$  is feasible for  $\min_{\nu \in \mathcal{M}(q)} \|\hat{\nu} - \hat{\mu} - \hat{\rho}\|_2$ . For the scaling of  $c_{d,\tau}^{(5)}$  we note at first that by the separation condition the maximal number of nodes  $M$  is bounded by  $M \leq \left(\frac{n}{q_\tau}\right)^d$  and thus  $c_{d,\tau}^{(5)} := c_{d,\tau}^{(4)} q_\tau^{-d/2}$ . From the proof of Theorem 2.2.13, we can derive

$$c_{d,\tau}^{(4)} \leq \frac{1}{4} \sqrt{\frac{3d}{\min\left(c_{d,\tau}^{(1)}, c_{d,\tau}^{(3)}, \frac{dc_{d,\tau}^{(2)}q_\tau^2}{32q^2}\right)}} \leq \frac{1}{4} \sqrt{\frac{3d}{\min\left(c_{d,\tau}^{(1)}, c_{d,\tau}^{(3)}, \frac{dc_{d,\tau}^{(2)}q_\tau^2}{8}\right)}} \lesssim_d \tau^{-1/2} (1+\tau)^{d/4+1}.$$

Together with  $q_\tau^{-d/2} = (1+\tau)^{-d/4} \left(\frac{\pi}{j_{d/2,1}}\right)^{d/2}$  this gives the proposed order of  $c_{d,\tau}^{(5)}$ .  $\square$

By the previous corollary we see that the condition number stays bounded if the separation parameter  $q$  fulfils  $q > \frac{j_{d/2,1}}{\pi n}$ . On the contrary, this does not remain valid if  $q$  is slightly smaller as we will present in Theorem 2.2.21. Therefore, it is natural to define a *diffraction limit* by means of the condition number that distinguishes the two cases of a polynomial condition number and larger exponentially growing condition number for more densely spaced nodes. While there are many definitions for a diffraction limit in the special case  $d = 2$  motivated from applications in optics (see the introduction of this dissertation and [28] for an overview), this idea gives a mathematically rigid formulation. As it seems to be widely accepted that the diffraction limit is anti proportional to the cut off frequency  $n$ , we are interested to determine the constant depending on the dimension  $d$  where the condition is no longer a polynomial. In other words, we set  $q = \frac{\tilde{q}}{n}$  and study the optimal constant  $\tilde{q}$  as  $n \rightarrow \infty$ .

**Definition 2.2.18** (Diffraction limit). For  $n, d \in \mathbb{N}$  we define the *optimal transition constant*  $\Omega_d \geq 0$  as

$$\Omega_d = \inf \left\{ \tilde{q} > 0 : \exists \beta \in \mathbb{N} \lim_{n \rightarrow \infty} \sup_{M \leq (\sqrt{dn}/\tilde{q})^d} \frac{\kappa_{\text{abs}}\left(\frac{\tilde{q}}{n}, n, d, M\right)}{M^\beta} < \infty \right\}.$$

The *diffraction limit* for finite  $n$  is set to be  $\Omega_d n^{-1}$ .<sup>39</sup>

<sup>39</sup>Even though our aim was to present the diffraction limit as the transition between polynomial and exponential growth of the condition, this transition can only be seen as  $n \rightarrow \infty$  and thus we can just analyse the optimal transition constant  $\Omega_d$ . For finite  $n$  the ‘‘transition’’ heuristically happens at  $\Omega_d n^{-1} + o(n^{-1})$  and our definition of the diffraction limit omits the  $o$ -term.

We remark that the upper bound on the size of considered node sets  $M$  simply follows from the packing argument that one can pack less  $\ell^2$ -balls with radius  $q/2$  than  $\ell^\infty$ -balls with radius  $d^{-1/2}q/2$  into the torus  $\mathbb{T}^d$  while the latter problem can be solved easily. In order to find a lower bound on the transition constant  $\Omega_d$  for  $d = 2$ , we use the following construction by Chen and Moitra, see [28, Lem. 2.1]. In their paper, the diffraction limit was already defined similar to Definition 2.2.18. Unfortunately, we believe that their proof is not fully complete such that we add our own proof highlighting the previously missing details. But before we fix these issues, we use their method in order to obtain a lower bound on the transition constant in the univariate case. This lower bound was already found by Moitra for the sub-problem of estimating condition numbers of Vandermonde matrices, cf. [117, Thm. 3.1].

**Lemma 2.2.19** (Univariate lower bound). *Let  $\epsilon \in (0, 1)$ . Then, there exists a  $n_0 \in \mathbb{N}$  such that for all  $n \geq n_0$  there exist two univariate, nonnegative and  $n$ -sparse measures  $\mu_1, \mu_2 \in \mathcal{M}(q)$  where the separation  $q$  satisfies  $nq = 1 - \epsilon$  and  $W_1(\mu_1, \mu_2) \geq \frac{q}{2}$  while*

$$|\hat{\mu}_1(v) - \hat{\mu}_2(v)| \leq 8 \cdot \left( \frac{1}{\sqrt{2}} \right)^{\epsilon(n-1)}$$

for all  $v \in \mathbb{Z}$  with  $|v| \leq n$ .

*Proof.* Our proof technique is based on the proof of [28, Lem. 2.1]. We distinguish two cases depending on the parity of  $n$ . If  $n \in \mathbb{N}$  is even, we set  $x_j = \frac{jq}{2} = \frac{j(1-\epsilon)}{2n}$  for  $j = -n+1, \dots, n-1$ . Let  $F_n$  be the univariate Fejér kernel introduced in Definition 1.3.5,  $\lfloor \cdot \rfloor$  the floor function and set

$$H(v) = \frac{1}{(l+1)^{\lfloor (n-1)/l \rfloor}} F_l^{\lfloor (n-1)/l \rfloor} \left( \frac{(1-\epsilon)v}{2n} + \frac{1}{2} \right)$$

for some odd  $l \in \mathbb{N}$ ,  $l \leq n-1$ . By this definition,  $1 \leq \lfloor \frac{n-1}{l} \rfloor \in \mathbb{N}$  and thus  $H$  is a trigonometric polynomial with degree at most  $n-1$ . As

$$(n+1)^{-1} F_n(x) = \sum_{k=-n}^n \frac{n+1-|k|}{(n+1)^2} e^{2\pi i k x}$$

is a trigonometric polynomial with nonnegative coefficients summing to one, one finds by the relation between multiplication of functions and the convolution of their Fourier coefficients

$$\frac{1}{(l+1)^{\lfloor (n-1)/l \rfloor}} F_l^{\lfloor (n-1)/l \rfloor}(x) = \sum_{k=-n+1}^{n-1} \alpha_k e^{2\pi i k x}$$

for some  $\alpha_k \geq 0$  with  $\sum_k \alpha_k = 1$ . Inserting this into the definition of  $H$  gives

$$H(v) = \sum_{k=-n+1}^{n-1} \alpha_k (-1)^k e^{2\pi i (1-\epsilon)kv/(2n)} = \sum_{k=-n+1}^{n-1} \alpha_k (-1)^k e^{2\pi i v x_k}. \quad (2.24)$$

Because  $n$  is even, we represent

$$H(v) = \sum_{j=-n/2+1}^{n/2-1} \alpha_{2j} e^{2\pi i v x_{2j}} - \sum_{j=-n/2}^{n/2-1} \alpha_{2j+1} e^{2\pi i v x_{2j+1}} = \hat{\nu}_1(v) - \hat{\nu}_2(v) \quad (2.25)$$

## 2 Condition of sparse super resolution

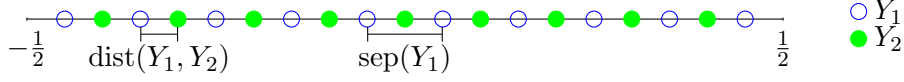


Figure 2.7: Definition of alternating node sets  $Y_1, Y_2$  for  $n = 10$ . As  $n$  is even,  $|Y_1| = n$  and  $|Y_2| = n - 1$ . The cardinality of  $Y_1$  and  $Y_2$  is interchanged for odd  $n$ .

as the difference of moments of two measures  $\nu_1, \nu_2$  supported on  $Y_1 = \{x_j : |j| \text{ even}, |j| \leq n - 1\}$ ,  $Y_2 = \{x_j : |j| \text{ odd}, |j| \leq n - 1\}$ . We display  $Y_1$  and  $Y_2$  in Figure 2.7. Since the nonnegative weights  $\alpha_k$  sum to one and satisfy

$$\sum_{j=-n/2+1}^{n/2-1} \alpha_{2j} - \sum_{j=-n/2}^{n/2-1} \alpha_{2j+1} = H(0) = \frac{\sin((l+1)\frac{1}{2}\pi)^{2\lfloor n/l \rfloor}}{(l+1)^{\lfloor n/l \rfloor} \sin(\frac{1}{2}\pi)^{2\lfloor n/l \rfloor}} = 0$$

due to  $l$  being odd, we can multiply each weight of  $\nu_1, \nu_2$  with

$$\left( \sum_{j=-n/2+1}^{n/2-1} \alpha_{2j} \right)^{-1} = \left( \sum_{j=-n/2}^{n/2-1} \alpha_{2j+1} \right)^{-1} = 2$$

in order to obtain nonnegative measures  $\mu_1, \mu_2$  with coefficients summing to one. Each of the measures has separation  $q = (1 - \epsilon)n^{-1}$  by construction and up to  $M = n$  nodes. Furthermore, all nodes of  $\mu_1$  and  $\mu_2$  are separated by at least  $\frac{q}{2}$  leading to the proposed lower bound on their Wasserstein distance. For all  $v$  with  $-n \leq v \leq n$  we can derive  $\frac{\epsilon}{2} \leq \frac{(1-\epsilon)v}{2n} + \frac{1}{2} \leq 1 - \frac{\epsilon}{2}$  giving  $\|\frac{(1-\epsilon)v}{2n} + \frac{1}{2}\|_{\mathbb{T}} \geq \frac{\epsilon}{2}$ . It is well known that

$$|F_n(x)| \leq \frac{1}{n+1} \frac{1}{4\|x\|_{\mathbb{T}}^2}$$

by estimating the sine function, e.g. see [132, p. 25]. Therefore, we end up with

$$|\hat{\mu}_1(v) - \hat{\mu}_2(v)| = 2|\hat{\nu}_1(v) - \hat{\nu}_2(v)| = 2H(v) \leq \frac{2}{((l+1)\epsilon)^{2\lfloor (n-1)/l \rfloor}}.$$

For given  $\epsilon > 0$  we choose  $n'_0$  as the smallest integer such that  $n'_0\epsilon \geq 2$ . Setting

$$l := \begin{cases} n'_0 + 1, & n'_0 \text{ even,} \\ n'_0, & n'_0 \text{ odd} \end{cases} \quad (2.26)$$

gives an odd  $l$  as desired. Moreover, we have  $l + 1 \leq \frac{2}{\epsilon} + 2$  and  $l \leq n'_0 + 1 \leq n - 1$  for all even  $n \geq n_0 := n'_0 + 2$ . This gives then

$$|\hat{\mu}_1(v) - \hat{\mu}_2(v)| \leq 2 \cdot 2^{-2\lfloor \frac{n-1}{l} \rfloor} \leq 8 \cdot 2^{-\frac{2(n-1)}{2/\epsilon+2}} = 8 \cdot 2^{-\epsilon(n-1)(1+\epsilon)^{-1}} \leq 8 \cdot 2^{-\frac{1}{2}\epsilon(n-1)}$$

by using  $\epsilon < 1$  in the last step. For the other case of  $n$  being odd, one just arranges the terms for  $\nu_1, \nu_2$  in (2.25) slightly different according to their sign whereas all other calculations remain valid independent of the parity of  $n$ .  $\square$

Analogously to this lemma, we can refine [28, Lem. 2.1].

**Lemma 2.2.20.** (Lower bound by Chen and Moitra, cf. [28, Lem. 2.1]) Let  $\epsilon \in (0, 1)$ . Then, there exists a  $n_0 \in \mathbb{N}$  such that for all  $n \geq n_0$  there exist two bivariate, nonnegative and  $2n^2$ -sparse measures  $\mu_1, \mu_2 \in \mathcal{M}(q)$  where the separation  $q$  satisfies  $nq = \sqrt{\frac{4}{3}}(1 - \epsilon)$  and  $W_1(\mu_1, \mu_2) \geq \frac{q}{2}$  while

$$|\hat{\mu}_1(k) - \hat{\mu}_2(k)| \leq 8 \cdot 2^{-\frac{1}{2}\epsilon(n-1)}$$

for all  $k \in \mathbb{Z}^2$  with  $\|k\|_2 \leq n$ .

*Proof.* Again, two measures with small distance of their moments sequence need to be constructed. In this bivariate setting, we generalise the equidistant nodes from the univariate example, see the proof of Lemma 2.2.19, to nodes on a lattice such that each of the two measures is supported on a hexagonal lattice. In our notation, we define

$$H(v_1, v_2) = \frac{F_l^{\lfloor (n-1)/l \rfloor} \left( \frac{(1-\epsilon)v_1}{2n} + \frac{1}{2} \right) F_l^{\lfloor (n-1)/(\sqrt{3}l) \rfloor} \left( \frac{(1-\epsilon)\sqrt{3}v_2}{2n} + \frac{1}{2} \right)}{(l+1)^{\lfloor (n-1)/l \rfloor + \lfloor (n-1)/(\sqrt{3}l) \rfloor}}$$

with the same choice of  $n'_0, l, n_0$  and  $n \geq n_0$  as in the proof of Lemma 2.2.19, cf. (2.26).<sup>40</sup> Inserting the polynomial expansion as in (2.24), we derive

$$H(v_1, v_2) = \sum_{k_1=-n+1}^{n-1} \sum_{k_2=-\lfloor (n-1)/\sqrt{3} \rfloor}^{\lfloor (n-1)/\sqrt{3} \rfloor} \alpha_{k_1} \alpha_{k_2} (-1)^{k_1+k_2} e^{2\pi i \frac{(1-\epsilon)}{2n} (v_1 k_1 + \sqrt{3} v_2 k_2)}.$$

If we define a set of lattice points

$$Y := \left\{ \frac{(1-\epsilon)}{2n} \left( k_1, \sqrt{3}k_2 \right)^\top, k_1 = -n+1, \dots, n-1 \text{ and } k_2 = -\left\lfloor \frac{n-1}{\sqrt{3}} \right\rfloor, \dots, \left\lfloor \frac{n-1}{\sqrt{3}} \right\rfloor \right\},$$

one can rewrite  $H$  as the difference of moments corresponding to nonnegative measures  $\nu_1, \nu_2$  supported either on the elements of  $Y$  where  $k_1 + k_2$  is even or odd respectively. We display this in Figure 2.8 highlighting that the separation of each of the measures is equal to  $q$ . Again, the weights can be multiplied by two in order to generate probability measures  $\mu_1$  and  $\mu_2$ . Moreover, the lattice points of the different lattices are separated by at least  $\frac{q}{2}$ . Hence, we can easily bound  $W_1(\mu_1, \mu_2) \geq \frac{q}{2}$ . Each of the two measures is by definition supported on up to  $M = \left\lceil \frac{1}{2}(2n-1) \left( 2 \left\lfloor \frac{n-1}{\sqrt{3}} \right\rfloor + 1 \right) \right\rceil \leq \sqrt{3}n^2$  nodes. Now, we need a lower bound on the norm of the argument of the Fejér kernels in  $H$  in order to estimate this  $H$ . Even if there are inconsistent statements,<sup>41</sup> it was already discussed in [28, Lem. 2.1] that the set  $\left\{ \left( \frac{(1-\epsilon)v_1}{2n} + \frac{1}{2}, \frac{(1-\epsilon)\sqrt{3}v_2}{2n} + \frac{1}{2} \right)^\top : v \in \mathbb{R}^2 \text{ with } \|v\|_2 \leq n \right\}$  is an ellipsoid with axis lengths  $1 - \epsilon$  and  $\sqrt{3}(1 - \epsilon)$  around  $(-\frac{1}{2}, -\frac{1}{2})$ . By the symmetry of this ellipsoid and its length in the first coordinate, we have that no integer vector is contained

<sup>40</sup>We remark that the mentioned problem of [28, Lem. 2.1] arises at this point. They take powers of the Fejér kernel which are neither an integer nor larger than one. Hence, one cannot conclude by the convolution theorem that  $H$  is a polynomial. We omit this by taking the floor function  $\lfloor \cdot \rfloor$  and  $l \leq n-1$ . As we wanted to choose  $(l+1)\epsilon \geq 2$  in the later course of the proof, we need  $n$  to be large enough in contrast to  $\epsilon$ . So the statement of the theorem can be proven only for  $n \geq n_0$  for some  $n_0 \in \mathbb{N}$ . The latter condition was not included in [28, Lem. 2.1].

<sup>41</sup>The lower bound is claimed to be “ $\epsilon/2\sqrt{2}$ ” at first (cf. [28, p. 494]), while it is later said to be  $\epsilon/2$ , see [28, p. 495].

## 2 Condition of sparse super resolution

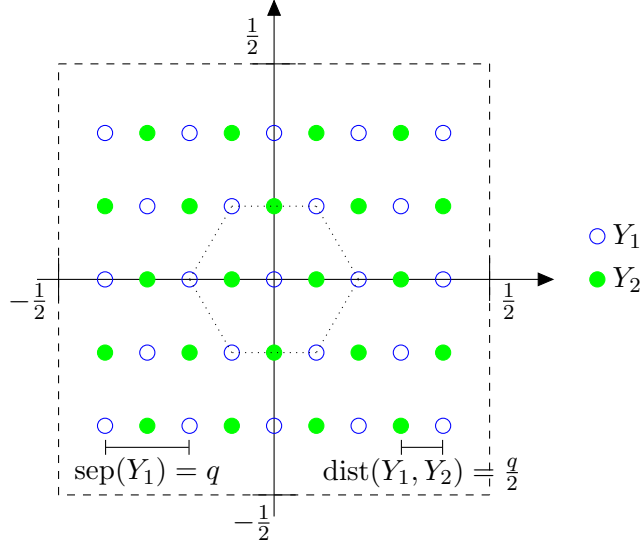


Figure 2.8: Support of the measures  $\nu_1$  and  $\nu_2$  for  $n = 5$ . We remark that both sets  $Y_1$  and  $Y_2$  are regular hexagonal lattices.

in it. Due to periodicity, we can study the distance to the origin and look for the element of this ellipsoid with the smallest  $\ell^\infty$ -norm. Necessarily, this element is on the boundary and parameterising this by some angle  $\phi \in [0, 2\pi)$  gives

$$\begin{aligned}
 & \min_{\phi} \min_{m \in \mathbb{Z}^2} \left\| \left( \frac{(1-\epsilon) \cos \phi}{2} + \frac{1}{2}, \frac{(1-\epsilon)\sqrt{3} \sin \phi}{2} + \frac{1}{2} \right)^\top + m \right\|_{\infty} \\
 &= \min_{\phi} \max \left( \frac{(1-\epsilon) \cos \phi}{2} + \frac{1}{2}, \left| \frac{(1-\epsilon)\sqrt{3} \sin \phi}{2} + \frac{1}{2} \right| \right) \\
 &\geq \min_{\phi} \frac{(1-\epsilon) \cos \phi}{2} + \frac{1}{2} \\
 &= \frac{\epsilon}{2}.
 \end{aligned}$$

Finally, this then allows us to bound

$$|\hat{\mu}_1(v) - \hat{\mu}_2(v)| = 2|\hat{\nu}_1(v) - \hat{\nu}_2(v)| = 2H(v) \leq \frac{2}{((l+1)\epsilon)^{2\lfloor (n-1)/l \rfloor}}$$

for all  $v \in \mathbb{R}^2$  with  $\|v\|_2 \leq n$ . As we have chosen  $l$  and  $n_0$  as in the proof of Lemma 2.2.19, we get the same estimate.  $\square$

In [28], the upper bound for the bivariate diffraction limit is  $n\Omega_2 \leq \frac{2j_{0,1}}{\pi} \approx 1.53$ . By Corollary 2.2.17, we improve this in the following theorem showing almost matching upper and lower bounds for  $d = 2$ . The simpler univariate case can be analysed completely.

**Theorem 2.2.21** (Bounds on the diffraction limit). *In any dimension  $d \in \mathbb{N}$ , the optimal transition constant satisfies*

$$\Omega_d \leq \frac{j_{d/2,1}}{\pi}.$$



This bound is sharp in the univariate case, i.e.  $\Omega_1 = 1$  for  $d = 1$ , and the bivariate diffraction limit relevant in imaging applications can be estimated by

$$1.16 \approx \sqrt{\frac{4}{3}} \leq \Omega_2 \leq \frac{j_{1,1}}{\pi} \approx 1.22.$$

*Proof.* For  $nq > \frac{j_{d/2,1}}{\pi}$  we have shown in Corollary 2.2.17 that

$$\kappa_{\text{abs}}(q, n, d, M) \lesssim_d \frac{2(1+\tau)}{\tau} \leq \frac{2 \left(\frac{\pi}{j_{d/2,1}}\right)^2 n^2 q^2}{\left(\frac{\pi}{j_{d/2,1}}\right)^2 n^2 q^2 - 1} = 2 + \frac{2}{\left(\frac{\pi}{j_{d/2,1}}\right)^2 n^2 q^2 - 1}.$$

In other words, we can bound  $\kappa_{\text{abs}}(q, n, d, M)$  for  $nq > \frac{j_{d/2,1}}{\pi}$  polynomially in  $(nq)^{-1}$  and thus we can take  $\beta = 0$  in the definition of  $\Omega_d$ . Then, the infimum gives the first part of the theorem. The lower bounds for  $d = 1, 2$  can be derived from Lemma 2.2.19 or Lemma 2.2.20 respectively because for  $n$  large enough we explicitly constructed  $\mu_1, \mu_2$  such that

$$\sup_{M \leq (\sqrt{dn}/\tilde{q})^d} \frac{\kappa_{\text{abs}}\left(\frac{\tilde{q}}{n}, n, d, M\right)}{M^\beta} \geq \frac{W_1(\mu_1, \mu_2)}{\|\hat{\mu}_1 - \hat{\mu}_2\|_2 (\sqrt{dn}/\tilde{q})^{d\beta}} \geq \frac{(1-\epsilon)2^{\frac{1}{2}\epsilon(n-1)}}{16(\sqrt{dn}/\tilde{q})^{d\beta}}$$

and the lower bound goes to infinity as  $n \rightarrow \infty$  for any  $\beta \in \mathbb{N}$ . Making  $\epsilon$  arbitrarily small yields equality in the univariate situation and the almost sharp estimate for  $d = 2$ .<sup>42</sup>  $\square$

As remarked in the introduction, the bivariate upper bound  $\frac{j_{1,1}}{\pi} \approx 1.22$  on  $\Omega_d$  is known as *Rayleigh's criterion* in optics. Therefore, Theorem 2.2.21 gives a strong evidence that this criterion is also mathematically meaningful.<sup>43</sup> Another reason to promote the Rayleigh criterion is the statistical approach given by the *Cramér-Rao (CR) lower bound*. As described in Section 2.1, the CR bound estimates that the covariance of each unbiased estimator  $\hat{\theta}$  for a vector of parameters  $\theta$  can be bounded from below by the inverse of the *Fisher information matrix*  $J(\theta)$ . We summarise known results about the CR lower bound and the Fisher information matrix in the following theorem.

**Theorem 2.2.22.** (*CR and Fisher information, cf. [126, p. 6424]*) Assume that a random vector  $y \in \mathbb{C}^m$  has probability density function  $f(y, \theta)$  depending on some unknown, deterministic parameter  $\theta \in \mathbb{R}^l$  for some  $l \in \mathbb{N}$ . Then, the Fisher information matrix defined as the covariance<sup>44</sup>

$$J(\theta) = \mathbb{E}_y \left[ \left( \frac{\partial \log f(y, \theta)}{\partial \theta} \right) \left( \frac{\partial \log f(y, \theta)}{\partial \theta} \right)^* \right] \in \mathbb{C}^{l \times l}$$

satisfies  $\mathbb{E}_y \left[ (\hat{\theta}(y) - \theta)(\hat{\theta}(y) - \theta)^* \right] \succeq J(\theta)^{-1}$  for any unbiased estimator  $\hat{\theta}$ . If  $y$  follows a multivariate complex normal distribution with mean  $x(\theta) \in \mathbb{C}^m$  and diagonal covariance

<sup>42</sup>One might consider to obtain lower bounds for  $d > 2$  by the idea of Lemma 2.2.20 or Lemma 2.2.19 respectively but on one hand this might not be so interesting from the applied point of view while on the other hand packing arguments needed for the generalisation to higher dimensions are not straightforward. Hence, it would not be realistic to hope for similarly sharp bounds.

<sup>43</sup>As there is at least with our proof technique no hope to reduce the upper bound below  $\frac{j_{d/2,1}}{\pi}$  (see Remark 2.2.10), one could conjecture  $\Omega_2 = \frac{j_{d/2,1}}{\pi}$ . In order to prove this, one would need an example of two less densely spaced measures with exponentially small  $\ell^2$  distance of their moments.

<sup>44</sup>We emphasise that the expectation is computed with respect to the random  $y$  by the subscript  $\mathbb{E}_y$ .

## 2 Condition of sparse super resolution

matrix  $\delta^2 I$  for some  $\delta > 0$ , i.e.  $y \sim \mathcal{CN}(x(\theta), \delta^2 I)$ , then the Fisher information matrix can be calculated as

$$J(\theta) = \delta^{-2} G^* G, \quad \text{where } G = \left[ \frac{\partial x(\theta)}{\partial \theta_1}, \dots, \frac{\partial x(\theta)}{\partial \theta_l} \right] \in \mathbb{C}^{m \times l}.$$

From a theoretical point of view, this theorem allows to derive a minimal covariance of any algorithm recovering the  $\theta$  from measurements  $y$ . Hence, this can be seen as a lower bound on the condition of the problem itself. This theory was therefore used in the context of univariate super resolution in [53] by assuming that the measured moments are

$$\hat{\mu}(k) = \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k} + \hat{\rho}(k) \quad (2.27)$$

for some normally distributed noise  $\hat{\rho} \sim \mathcal{CN}(0, \delta^2 I)$  and this can be directly generalised from the univariate case  $k \in \{-n, \dots, n\}$  to the higher dimensional case  $\mathcal{I} := \{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}$ . Here, the vector of unknown parameters  $\theta$  is

$$\theta = (\alpha, Y)^\top := [(\alpha_t)_{t \in Y}, (t_1)_{t \in Y} \cdots (t_d)_{t \in Y}]^\top \in \mathbb{C}^{|Y|(d+1)}$$

and the measurements are  $y = (\hat{\mu}(k))_{\{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}} \in \mathbb{C}^{|\mathcal{I}|}$ . Based on this model, we can compute the Fisher information matrix as follows.

**Corollary 2.2.23.** *If the moments satisfy the noise model (2.27), we have the factorisation  $J(\alpha, Y) = \delta^{-2} G^* G$  of the Fisher information matrix where*

$$G = \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) D_\alpha$$

with a Vandermonde matrix

$$\mathcal{A} = \left( e^{-2\pi i t k} \right)_{k \in \{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}, t \in Y} \in \mathbb{C}^{|\mathcal{I}| \times |Y|},$$

matrices  $\tilde{\mathcal{A}}_s$ ,  $s = 1, \dots, d$ , with

$$\tilde{\mathcal{A}}_s = -2\pi i \left( k_s e^{-2\pi i t k} \right)_{k \in \{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}, t \in Y} \in \mathbb{C}^{|\mathcal{I}| \times |Y|}$$

and the diagonal matrix

$$D_\alpha := \text{diag} \left( \underbrace{1, \dots, 1}_{|Y| \text{ times}}, \underbrace{\alpha_{t_1}, \dots, \alpha_{t_{|Y|}}, \dots, \alpha_{t_{|Y|}}, \dots, \alpha_{t_{|Y|}}}_{\text{repeat weight vector } d \text{ times}} \right) \in \mathbb{C}^{|Y|(d+1) \times |Y|(d+1)}.$$

The matrices  $\tilde{\mathcal{A}}_s$  can be seen as a variant of a confluent Vandermonde matrix (see [57]).

*Proof.* The univariate case  $d = 1$  was given in [53] and the higher dimensional result follows from Theorem 2.2.22 by differentiating (2.27) with respect to the parameters. The derivative with respect to the weights gives the Vandermonde matrix  $\mathcal{A}$  while the partial derivatives with respect to the  $s$ th component of every node gives the confluent Vandermonde matrix  $\tilde{\mathcal{A}}_s$ .  $\square$

As the inverse of Fisher information is then a lower bound for the covariance of each unbiased estimator, it is then natural to define the condition of super resolution through the size of  $J(\theta)^{-1}$  and the problem is considered to be ill-conditioned if  $\|J(\theta)^{-1}\|$  becomes large or equivalently  $\lambda_{\min}(J(\theta))$  is very small. Hence, one is interested to find lower bounds on the minimal eigenvalue of  $J(\theta)$  in order to establish well-conditionedness. This approach was introduced in [53] by defining the transition between good and ill-conditionedness as follows.<sup>45</sup>

**Definition 2.2.24.** (Condition via CR, generalisation of [53, Def. 1]) The super resolution problem is said to be well-conditioned for some  $\tilde{q} > 0$  if for all  $n$  and parameter configurations with separation  $n \cdot \text{sep}Y \geq \tilde{q}$  and some minimal absolute value of all weights  $\alpha_{\min} > 0$  there exists a constant  $c_{\tilde{q}, \alpha_{\min}}$  independent of  $n$  such that

$$n^{-d} \lambda_{\min}(J(\theta)) \geq \delta^{-2} c_{\tilde{q}, \alpha_{\min}}.$$

Due to Corollary 2.2.23, one directly finds

$$\lambda_{\min}(J(\theta)) = \delta^{-2} \sigma_{\min}^2(G) \geq \frac{\min(1, \alpha_{\min}^2)}{\delta^2} \sigma_{\min}^2(\mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d) \quad (2.28)$$

with equality if all weights are equal to one.<sup>46</sup> Consequently, the problems boils down to an estimate on the smallest singular value of a block matrix where each block consists of a Vandermonde matrix or of a confluent Vandermonde matrix. While there have been many attempts to analyse the smallest singular value of Vandermonde matrices, see [120] for an overview and Section 2.3 for a contribution in this work, Ferreira Da Costa and Mitra [53] observed already for the one dimensional case that this can be done by an admissible function as we defined it in Definition 2.2.1. Whereas they utilised a variant of the Beurling-Selberg minorant for this, we can apply the function  $\psi$  from Lemma 2.2.2 having optimally small support.

**Proposition 2.2.25.** (Conditioning of partially confluent block Vandermonde matrix, generalisation of [53, Prop. 6]) Assume that  $n$  and the separation  $q$  satisfy  $nq > \sqrt{1 + \tau} \frac{j_{d/2}}{\pi}$  for some  $\tau > 0$ . As in the proof of Theorem 2.2.8, we set  $\psi_{\tau, n}(x) := n^d \psi_{\tau}(n \cdot x)$  where  $\psi_{\tau}$  is the admissible function defined in Lemma 2.2.2. Then, we have

$$\sigma_{\min}^2(\mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d) \geq c_{d, \tau}^{(6)} n^d$$

for some constant  $c_{d, \tau}^{(6)} > 0$ .

The proof uses the following lemma.

**Lemma 2.2.26** (Evaluating derivatives at zero). Let  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a radial function, i.e.  $\psi(x) = h(\|x\|_2)$  for some univariate function  $h$ . Assume that  $\psi, h$  are twice continuously differentiable and that  $\psi$  is maximal in zero. Then, we have  $\left(\frac{\partial \psi_{\tau, n}}{\partial x_s}\right)(0) = 0$  for all  $s = 1, \dots, d$ . Moreover, one can find

$$d \cdot \left(\frac{\partial^2 \psi_{\tau, n}}{\partial x_s^2}\right)(0) = \Delta \psi(0) \quad \text{and} \quad \left(\frac{\partial^2 \psi_{\tau, n}}{\partial x_s \partial x_{s'}}\right)(0) = 0$$

for any  $s, s' \in \{1, \dots, d\}, s \neq s'$ .

<sup>45</sup>In [53], the authors use the term stability instead of condition. Since we distinguish between the two terms as explained in Section 1.1, we proceed by using the term condition.

<sup>46</sup>In particular, we remark at this point that this analysis again separates the dependency of the condition on the weights from the influence of the nodes.

## 2 Condition of sparse super resolution

*Proof.* The vanishing gradient follows directly from the extremum in zero. For the second derivatives one can calculate

$$\left( \frac{\partial^2 \psi_{\tau,n}}{\partial x_s \partial x_{s'}} \right) (x) = \frac{h'(\|x\|_2)}{\|x\|_2} \delta_{s,s'} + \left( h''(\|x\|_2) - \frac{h'(\|x\|_2)}{\|x\|_2} \right) \frac{x_s x_{s'}}{\|x\|_2^2}.$$

Because  $h'(\|x\|_2) = h'(0) + h''(\|x\|_2)\|x\|_2 + o(\|x\|_2)$  as  $\|x\|_2 \rightarrow 0$  and  $\frac{|x_s x_{s'}|}{\|x\|_2^2} \leq 1$ , the second part vanishes in zero. This yields that the mixed derivatives vanish in zero. Finally, the first term is independent of  $s$  if  $s = s'$ . This gives the remaining part of the statement.  $\square$

We can then return to the proof of Proposition 2.2.25.

*Proof of Proposition 2.2.25.* We follow the idea of the proof of [53, Prop.6]. By the variational representation from Lemma 1.1.6, we have to find a lower bound on the expression  $\|(\mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d) u\|_2$  for any normalised vector  $u$  with block structure  $u = (u_0^\top, u_1^\top, \dots, u_d^\top)^\top \in \mathbb{C}^{|Y|(d+1)}$  where  $u_s \in \mathbb{C}^{|Y|}$ ,  $s = 1, \dots, d$ . We set

$$\hat{\mu}_0(k) := \sum_{t \in Y} (u_0)_t e^{-2\pi i t k} \quad \text{and} \quad \hat{\mu}_s(k) := - \sum_{t \in Y} 2\pi i k_s (u_s)_t e^{-2\pi i t k}$$

for  $s = 1, \dots, d$  and  $k \in \mathcal{I}$ . Now we can compute

$$\begin{aligned} \hat{\psi}_{\tau,n}(0) \left\| \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) u \right\|_2^2 &= \hat{\psi}_{\tau,n}(0) \sum_{k \in \mathbb{Z}^d, \|k\|_2 \leq n} \left| \sum_{s=0}^d \hat{\mu}_s(k) \right|^2 \\ &\geq \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) \left| \sum_{s=0}^d \hat{\mu}_s(k) \right|^2 \\ &= \sum_{s,s'=0}^d \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) \hat{\mu}_s(k) \overline{\hat{\mu}_{s'}(k)} \\ &= S_1 + S_2 + S_3 + S_4 \end{aligned}$$

where the decomposition consists of

$$\begin{aligned} S_1 &= \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) |\hat{\mu}_0(k)|^2, \\ S_2 &= \sum_{s=1}^d 2\Re \left[ \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) \hat{\mu}_s(k) \overline{\hat{\mu}_0(k)} \right], \\ S_3 &= \sum_{s=1}^d \sum_{\substack{s'=1 \\ s' < s}}^d 2\Re \left[ \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) \hat{\mu}_s(k) \overline{\hat{\mu}_{s'}(k)} \right] \quad \text{and} \\ S_4 &= \sum_{s=1}^d \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) |\hat{\mu}_s(k)|^2. \end{aligned}$$

By Poisson's summation formula and the separation of  $Y$  together with the compact support of  $\psi_{\tau,n}$  we derive

$$S_1 = \sum_{t,t' \in Y} (u_0)_t \overline{(u_0)_{t'}} \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau,n}(k) e^{2\pi i(t'-t)k} = \sum_{t \in Y} |(u_0)_t|^2 \psi_{\tau,n}(0)$$

and analogously due to the relation between multiplication with monomials and derivatives under the Fourier transform<sup>47</sup>

$$\begin{aligned}
 S_4 &= - \sum_{s=1}^d \sum_{t, t' \in Y} (u_s)_t \overline{(u_s)_{t'}} \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau, n}(k) (2\pi i k_s)^2 e^{2\pi i (t' - t)k} \\
 &= - \sum_{s=1}^d \sum_{t, t' \in Y} (u_s)_t \overline{(u_s)_{t'}} \sum_{k \in \mathbb{Z}^d} \left( \frac{\partial^2 \psi_{\tau, n}}{\partial x_s^2} \right) (k) e^{2\pi i (t' - t)k} \\
 &= - \sum_{s=1}^d \sum_{t \in Y} |(u_s)_t|^2 \left( \frac{\partial^2 \psi_{\tau, n}}{\partial x_s^2} \right) (0).
 \end{aligned}$$

Moreover, one can evaluate the cross terms  $S_2$  and  $S_3$  by observing

$$\begin{aligned}
 \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau, n}(k) \hat{\mu}_s(k) \overline{\hat{\mu}_0(k)} &= \sum_{t, t'} (u_s)_t \overline{(u_0)_{t'}} \sum_{k \in \mathbb{Z}^d} (-2\pi i k_s) \hat{\psi}_{\tau, n}(k) e^{2\pi i (t' - t)k} \\
 &= \sum_t (u_s)_t \overline{(u_0)_t} \left( \frac{\partial \psi_{\tau, n}}{\partial x_s} \right) (0)
 \end{aligned}$$

for  $s = 1, \dots, d$  and

$$\begin{aligned}
 \sum_{k \in \mathbb{Z}^d} \hat{\psi}_{\tau, n}(k) \hat{\mu}_s(k) \overline{\hat{\mu}_{s'}(k)} &= \sum_{t, t'} (u_s)_t \overline{(u_{s'})_{t'}} \sum_{k \in \mathbb{Z}^d} (-2\pi i k_s) (2\pi i k_{s'}) \hat{\psi}_{\tau, n}(k) e^{2\pi i (t' - t)k} \\
 &= - \sum_t (u_s)_t \overline{(u_{s'})_t} \left( \frac{\partial^2 \psi_{\tau, n}}{\partial x_s \partial x_{s'}} \right) (0)
 \end{aligned}$$

for  $s, s' \in \{1, \dots, d\}, s \neq s'$ . By Lemma 2.2.26, we have  $S_2 = S_3 = 0$  and

$$\begin{aligned}
 \left\| \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) u \right\|_2^2 &\geq \min \left( \frac{\psi_{\tau, n}(0)}{\hat{\psi}_{\tau, n}(0)}, - \frac{\left( \frac{\partial^2 \psi_{\tau, n}}{\partial x_s^2} \right) (0)}{\hat{\psi}_{\tau, n}(0)} \right) \|u\|_2^2 \\
 &\geq \min \left( c_{d, \tau}^{(1)}, \frac{c_d \tau (1 + \tau)^{-d/2-1} n^2}{4\pi^2 (1 + \tau) \hat{\varphi}(0)^2} \right) n^d
 \end{aligned}$$

where  $c_{d, \tau}^{(1)}$  is the constant from the proof of Theorem 2.2.8. Defining the constant given by the minimum as  $c_{d, \tau}^{(6)}$  completes the proof.  $\square$

As a corollary of Proposition 2.2.25, we obtain another argument for using the Rayleigh limit  $\frac{j_{d/2, 1}}{\pi n}$  in order to describe the stability of super resolution.

**Corollary 2.2.27.** *Let  $d \in \mathbb{N}$ . For all  $\tilde{q} > \frac{j_{d/2, 1}}{\pi}$ , the super resolution problem is well conditioned in the sense of Definition 2.2.24.*

*Proof.* This follows directly from Definition 2.2.24, Proposition 2.2.25 and (2.28).  $\square$

<sup>47</sup>Note that we estimated  $\hat{\psi}_{\tau, n}(k) \in \mathcal{O}(\|k\|_2^{-d-3})$  in footnote 30. Hence, this function allows to apply Poisson summation formula even to its second derivative.

## 2 Condition of sparse super resolution

In the univariate case, the sufficient condition from Corollary 2.2.27 for well-conditionedness reads  $qn = \tilde{q} > 1$  and this was already conjectured in [53] where this conjecture was formulated in an asymptotically equivalent way as  $q(2n + 1) > 2$ . Moreover, [53, Fig. 2] gives at least numerical evidence that  $\tilde{q} > 1$  is also necessary in the univariate situation. An approach to make this more precise by estimating the smallest singular value of  $(\mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d)$  from above can be done by using results on  $\sigma_{\min}(\mathcal{A})$ . In fact, choosing  $u = (u_0^\top, 0, \dots, 0)^\top \in \mathbb{C}^{|Y|(d+1)}$  where  $u_0$  is the normalised right singular vector corresponding to the smallest singular value of  $\mathcal{A}$  and  $\alpha = (1, \dots, 1)^\top \in \mathbb{C}^{|Y|}$  leads to

$$\delta^2 \lambda_{\min}(J(\alpha, Y)) = \sigma_{\min}^2 \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) \leq \left\| \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) u \right\|_2^2 = \sigma_{\min}^2(\mathcal{A}). \quad (2.29)$$

Even if the smallest singular values of Vandermonde matrices are well-studied, e.g. see [120] and the references therein, upper bounds on the smallest singular value for the case of ill-separated nodes in higher dimensions are difficult in general (cf. [120, Subsec. 3.4.4]). Nevertheless, the analysis from (2.29) together with Lemma 2.2.19 and Lemma 2.2.20 shows that the super resolution problem cannot be well conditioned in the sense of Definition 2.2.24 for  $\tilde{q} < \frac{1}{2}$  and  $\tilde{q} < \frac{1}{\sqrt{3}}$  in  $d = 1$  or  $d = 2$  respectively. Consequently, the stochastic approach with the Cramer-Rao bound gives the same bounds as Theorem 2.2.21 for the point where the condition of super resolution transitions from well- to ill-conditionedness.

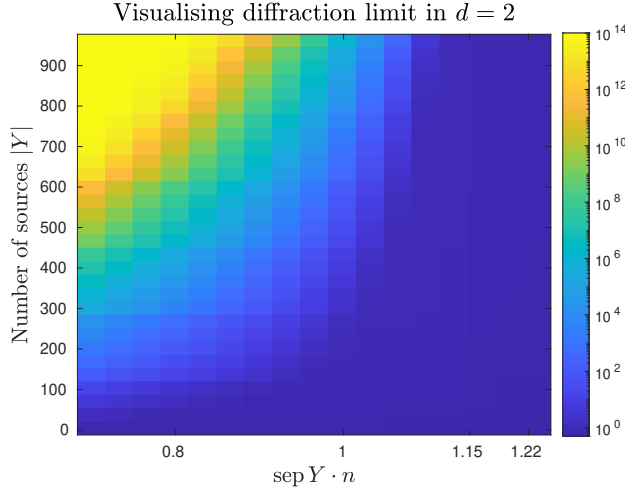


Figure 2.9: Visualisation of the bivariate diffraction limit. We place the node set  $Y$  on (a subset of) the two hexagonal lattices in  $\mathbb{T}^2$  from Lemma 2.2.20 while varying the separation  $\text{sep } Y$  and the number of nodes  $|Y|$  for fixed  $n = 40$ . For each selection of  $\text{sep } Y$  and  $|Y|$  we compute  $n \cdot \sigma_{\min} \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right)^{-1}$  as a proxy for the condition of the super resolution problem.

**Remark 2.2.28** (Condition for differentiable maps). We mentioned already in Chapter 1 that the absolute condition number of a differentiable map  $\phi$  is given by the norm of its Jacobian matrix  $\phi'$  (cf. [149, p. 90]). For our analysis of super resolution governed by the condition of the reconstruction map  $\mathcal{R}$  it might thus natural to study the moment map as the mapping of parameters  $t, \alpha_t, t \in Y$  to  $\hat{\mu}(k)$  and to compute the norm of its inverse by

the inverse function theorem. This would again motivate to define

$$\left\| \left( \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right) D_\alpha \right)^+ \right\|_2 \leq \max(1, \alpha_{\min}^{-1}) \cdot \sigma_{\min} \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right)^{-1}$$

as the condition of the reconstruction map  $\mathcal{R}$ . According to Corollary 2.2.27, this gives another justification for the Rayleigh limit. Furthermore, the formulation of the condition in terms of singular values of certain matrices allows to compute this condition for visualisation in Figure 2.9. In this numerical example, we see that the proxy  $n \cdot \sigma_{\min} \left( \mathcal{A}, \tilde{\mathcal{A}}_1, \dots, \tilde{\mathcal{A}}_d \right)^{-1}$  for the condition number of super resolution can become large if  $\text{sep} Y \cdot n < \frac{j_1+1}{\pi} \approx 1.22$ . However, the application of a variant of an inverse function theorem would require more justification such that we adhere to our original definition of a condition number given in Definition 2.2.16.

### 2.2.4 Condition of full inverse problem

In practice, one is typically not confronted with moments of a measure on the torus  $\mathbb{T}^d$  but with a perturbed low-pass version of a compactly supported measure which can be modelled as

$$g(x) = (h * \mu)(x) = \sum_{t \in Y} \alpha_t h(x - t), \quad x \in \mathbb{R}^d \quad (2.30)$$

with a PSF  $h \in L^1(\mathbb{R}^d) \cap C_0(\mathbb{R}^d) \cap \mathcal{B}_n(\mathbb{R}^d)$  for some  $n > 0$  and a discrete measure  $\mu = \sum_{t \in Y} \alpha_t \delta_t \in \mathcal{M} \left( [-\frac{1}{2}, \frac{1}{2}]^d \right)$ . The assumption of a bandlimited PSF is not too restrictive as many optical systems are usually bandlimited. Being interested in applications to microscopy images, we additionally restrict ourselves to the analysis of  $d = 2$  in this subsection and remark that a higher dimensional analogue could be derived similarly.

We have to incorporate into the model (2.30) that digital images consist of pixels and that at least in theory the values at the pixels are not exactly evaluations at a certain point  $x_j$  but an integral of  $g$  over the pixel  $x_j$ , i.e. a local mean of  $g$ . For example, this was mentioned in [75]. However, one could also circumvent this local mean by incorporating the effect of the integration into the PSF, see [75, p. 8]. Therefore we assume access to samples  $g(x_j) + \rho_j$  for  $x_j = \left( \frac{j_1}{2J}, \frac{j_2}{2J} \right)$ ,  $j \in \mathbb{Z}^2$ , such that  $x_j \in [-1/2 - \Delta, 1/2 + \Delta]^2$  originating from evaluation on some pixel grid in  $[-1/2 - \Delta, 1/2 + \Delta]^2$  with sampling parameter  $J \in \mathbb{N}$ , error  $\rho_j$  and field-of-view parameter  $\Delta > 0$ . Therein, the noise  $\rho_j$  is assumed to be bounded by a deterministic constant  $\varrho$  or to be stochastic while a sufficiently large field-of-view parameter  $\Delta$  is necessary in order to guarantee that no information on the boundary is lost. Then, the *full inverse problem of optical super resolution* would be to recover  $\mu$  from the measurements

$$\tilde{g} := (g(x_j) + \rho_j)_{x_j = \frac{j}{2J} \in [-1/2 - \Delta, 1/2 + \Delta]^2}$$

and we now want study the condition of this problem based on our findings from the previous subsection. In contrast to the previous analysis, we do not assume a periodic setting but are interested in results on the compact domain  $[-\frac{1}{2}, \frac{1}{2}]^2$  as a subset of  $\mathbb{R}^2$ . The main question of this subsection is then whether this changes the condition of the problem and the notion of our definition for a diffraction limit. We will see that this is not the case and that the main results of the previous subsection carry over to the compact domain setting. For the sake of simplicity, we make the following assumptions on the bandlimited PSF  $h$ .

## 2 Condition of sparse super resolution

**Definition 2.2.29.** We consider the PSF  $h \in L^1(\mathbb{R}^2) \cap C_0(\mathbb{R}^2) \cap \mathcal{B}_n(\mathbb{R}^2)$

1. to be a radial function with decay

$$0 \leq h(x) \leq \frac{c_1 n^2}{(1 + n\|x\|_2)^3} \quad \text{for some } c_1 > 0, \quad (2.31)$$

2. and to have radial derivative satisfying

$$\left| \frac{\partial h}{\partial r}(x) \right| \leq \frac{c_2 n^3}{(1 + n\|x\|_2)^3} \quad \text{for some } c_2 > 0, \quad (2.32)$$

3. while its Fourier transform admits

$$1 \geq \max_{\|v\| \leq n} |\hat{h}(v)|^2 > \min_{\|v\| \leq n'} |\hat{h}(v)|^2 \geq \frac{c_3 (n - n')^4}{n^4} \quad (2.33)$$

for some  $c_3 > 0$  and any  $0 < n' < n$ .

We count the number of pixels by  $\mathcal{J} := \left\{ j \in \mathbb{Z}^2 : x_j = \frac{j}{2^J} \in [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2 \right\}$ . In analogy to Definition 2.2.15 and Definition 2.2.16 we define the following.

**Definition 2.2.30** (Reconstruction from image data). For given PSF  $h$  as in Definition 2.2.29 and  $q, \Delta > 0$  the *image data reconstruction map* is

$$\tilde{\mathcal{R}} : \mathbb{C}^{\mathcal{J}} \rightarrow \mathcal{P}(\mathcal{M}(q)), \quad \tilde{g} \mapsto \operatorname{argmin}_{\nu \in \mathcal{M}(q)} \sum_{j \in \{0, \dots, J-1\}^d} |\tilde{g}_j - (h * \nu)(x_j)|^2$$

where we consider in this section  $\mathcal{M}(q)$  as the set of *non-periodic*, discrete measures with support  $Y \subset [-\frac{1}{2}, \frac{1}{2}]^2$  having *Euclidean* separation  $\min_{t, t' \in Y} \|t - t'\|_2$  at least  $q$ .

**Definition 2.2.31** (Condition for full inverse problem of SR). We define the *condition number of recovery from image data* as<sup>48</sup>

$$\tilde{\kappa}_{\text{abs}}(q, \Delta, J, h, M) := \sup_{\substack{\mu \in \mathcal{M}(q) \\ |Y^\mu| \leq M}} \sup_{\substack{\rho \in \mathbb{C}^{\mathcal{J}} \\ \rho \neq 0}} \inf_{\nu \in \mathcal{R}(((h * \mu)(x_j))_{j+\rho})} \frac{W_1(\nu, \mu)}{\|\rho\|_2}.$$

After defining the condition number for this type of problem, we are of course interested to relate it to our results from the previous subsection in order to understand whether the problem is different if we take a PSF-convolved image instead of truncated moments into account. Similar to the previous approach, we then want to control the difference  $\|\tilde{g}^{(1)} - \tilde{g}^{(2)}\|_2$  where  $\tilde{g}_j^{(1)} = (h * \mu_1)(x_j)$  for some  $\mu_1 \in \mathcal{M}(q)$  and  $\tilde{g}^{(2)}$  is analogously constructed from some measure  $\mu_2 \in \mathcal{M}(q)$ . Hence, we calculate

$$\begin{aligned} \frac{\|\tilde{g}^{(1)} - \tilde{g}^{(2)}\|_2^2}{\sum_{j' \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_{j'})|^2} &= \frac{\sum_{x_j \in [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} |(h * \mu_1)(x_j) - (h * \mu_2)(x_j)|^2}{\sum_{j' \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_{j'})|^2} \\ &= 1 - \frac{\sum_{x_j \notin [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} |(h * (\mu_1 - \mu_2))(x_j)|^2}{\sum_{j \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_j)|^2} \end{aligned} \quad (2.34)$$

<sup>48</sup>In this section, we use the 1-Wasserstein distance according to Proposition 1.4.5 with  $\mathcal{X} = [-\frac{1}{2}, \frac{1}{2}]^2$  equipped with the Euclidean distance.



where the sum over  $\mathbb{Z}^d$  can be estimated by Lemma 1.2.5. This gives then

$$\begin{aligned} \sum_{j \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_j)|^2 &= 4J^2 \int_{\mathbb{R}^2} |(h * (\mu_1 - \mu_2))(x)|^2 dx \\ &= 4J^2 \int_{\mathbb{R}^2} |\hat{h}(v)|^2 |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv \\ &\geq 4J^2 \left( \min_{\|v\| \leq n'} |\hat{h}(v)|^2 \right) \int_{B_{n'}(0)} |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv \end{aligned}$$

if  $J \geq n > n'$  for some  $n' \in \mathbb{R}$ . Thus, control over the denominator in (2.34) is possible by using (2.33) and formulating an analogue to Theorem 2.2.8 for the low order  $L^2$ -difference of the Fourier transforms of  $\mu_1$  and  $\mu_2$ .

**Lemma 2.2.32** (Non-periodic Ingham inequality). *Let  $J \geq n > n'$  for a spatial sampling parameter  $J \in \mathbb{N}$  and  $n, n' > 0$ . For  $\mu_1, \mu_2 \in \mathcal{M}(q)$  with*

$$q = \frac{\sqrt{1 + \tau} j_{1,1}}{\pi n'} = \frac{\sqrt{1 + \tau} j_{1,1}}{\pi \gamma n} \quad (2.35)$$

for some  $\tau > 0$  and  $\gamma \in (0, 1)$  one defines a disjoint decomposition of  $Y := \text{supp}(\mu_1 - \mu_2) = Y^{\mu_1} \cup Y^{\mu_2}$  into  $Y_1 \subset Y^{\mu_1}$ ,  $Y_2 \subset Y^{\mu_2}$  and  $Y_3 \subset Y^{\mu_1} \cup Y^{\mu_2}$ , see also [38, Thm. 3.6] and the proof of Theorem 2.2.8, with:

- (i)  $Y_3 := \{t \in Y : \text{For all } t' \in Y \text{ with } t \neq t' \text{ one has } \|t - t'\|_2 \geq \frac{q}{2}\}$
- (ii) For all  $t \in Y_1$  there is exactly one  $\eta(t) \in Y_2$  with  $\|t - \eta(t)\|_2 < \frac{q}{2}$ .

Then, we can use the same constants  $c_{d,\tau}^{(1)}, c_{d,\tau}^{(2)}, c_{d,\tau}^{(3)} > 0$  and the notation  $\tilde{\alpha}_t$  from Theorem 2.2.8 for the estimate

$$\begin{aligned} \int_{B_{n'}(0)} |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv &\geq \sum_{t \in Y_3} c_{d,\tau}^{(1)} n'^d |\tilde{\alpha}_t|^2 \\ &+ \sum_{t \in Y_1} \frac{1}{2} c_{d,\tau}^{(2)} n'^{d+2} \left( |\alpha_t^{(1)}|^2 + |\alpha_t^{(2)}|^2 \right) \|t - \eta(t)\|_2^2 + c_{d,\tau}^{(3)} n'^d |\alpha_t^{(1)} - \alpha_{\eta(t)}^{(2)}|^2. \end{aligned}$$

*Proof.* Beginning with the computation

$$\begin{aligned} \int_{B_{n'}(0)} |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv &\geq \hat{\psi}_{\tau, n'}(0)^{-1} \int_{\mathbb{R}^2} \hat{\psi}_{\tau, n'}(v) |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv \\ &= \hat{\psi}_{\tau, n'}(0)^{-1} \sum_{t, t' \in Y} \tilde{\alpha}_t \overline{\tilde{\alpha}_{t'}} \psi_{\tau, n'}(t - t') \end{aligned}$$

for  $\psi_{\tau, n'}$  as in the proof of Theorem 2.2.8, we observe that the inverse Fourier transformation can be used analogously to the Poisson summation formula in the periodic case. Consequently, the argument of  $\psi_{\tau, n'}$  depends only on the Euclidean distance instead of the wrap-around-metric on the torus. The rest of the proof is completely analogous to the proof of Theorem 2.2.8.  $\square$

We now turn to control over the numerator in (2.34).

## 2 Condition of sparse super resolution

**Lemma 2.2.33.** *Let  $h$  satisfy the assumptions of Definition 2.2.29. Under the assumptions and with the notation of Lemma 2.2.32 one can bound*

$$|h * (\mu_1 - \mu_2)(x)| \leq \frac{\max(c_1, \sqrt{2}c_2)n^2(\sum_{t \in Y_3} |\tilde{\alpha}_t| + \sum_{t \in Y_1} |\alpha_t - \alpha_{\eta(t)}| + n|\alpha_{\eta(t)}| \|t - \eta(t)\|_2)}{(1 + n \operatorname{dist}(x, Y))^3}$$

for all  $x \in \mathbb{R}^2$  with  $\|x\|_\infty \geq \frac{1}{2} + \Delta$ . Additionally, we can estimate

$$\frac{\max(c_1, \sqrt{2}c_2)^{-2} \sum_{x_j \notin [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} |(h * (\mu_1 - \mu_2))(x_j)|^2}{\sum_{t \in Y_3} |\tilde{\alpha}_t|^2 + \sum_{t \in Y_1} |\alpha_t - \alpha_{\eta(t)}|^2 + n^2(|\alpha_t|^2 + |\alpha_{\eta(t)}|^2) \|t - \eta(t)\|_2^2} \leq \frac{6(1 + \frac{1}{2\Delta})J^2|Y|n^2}{(1 + \Delta n)^4}.$$

*Proof.* For the first inequality we compute

$$\begin{aligned} & |h * (\mu_1 - \mu_2)(x)| \\ &= \left| \sum_{t \in Y_3} \tilde{\alpha}_t h(x - t) + \sum_{t \in Y_1} (\alpha_t - \alpha_{\eta(t)}) h(x - t) + \alpha_{\eta(t)} (h(x - t) - h(x - \eta(t))) \right| \\ &\leq \frac{c_1 n^2}{(1 + n \operatorname{dist}(x, Y))^3} \left( \sum_{t \in Y_3} |\tilde{\alpha}_t| + \sum_{t \in Y_1} |\alpha_t - \alpha_{\eta(t)}| \right) \\ &\quad + \sum_{t \in Y_1} |\alpha_{\eta(t)}| \left| \left\langle \int_0^1 \left( \frac{\partial h}{\partial x_1}, \frac{\partial h}{\partial x_2} \right)^\top (x - t + s(\eta(t) - t)) ds, t - \eta(t) \right\rangle \right| \end{aligned}$$

through Taylor's theorem. The Cauchy-Schwarz inequality, the chain rule, the mean value theorem and the assumptions of Definition 2.2.29 allow then to bound the inner product by

$$\begin{aligned} & \left\| \int_0^1 \left( \frac{\partial h}{\partial x_1}, \frac{\partial h}{\partial x_2} \right)^\top (x - t + s(\eta(t) - t)) ds \right\|_2 \|t - \eta(t)\|_2 \\ &\leq \left\| \int_0^1 \frac{\partial h}{\partial r} (x - t + s(\eta(t) - t)) \frac{x - t + s(\eta(t) - t)}{\|x - t + s(\eta(t) - t)\|_2} ds \right\|_1 \|t - \eta(t)\|_2 \\ &\leq \int_0^1 \frac{\partial h}{\partial r} (x - t + s(\eta(t) - t)) \frac{\|x - t + s(\eta(t) - t)\|_1}{\|x - t + s(\eta(t) - t)\|_2} ds \cdot \|t - \eta(t)\|_2 \\ &\leq \frac{\sqrt{2}c_2 n^3}{(1 + n \operatorname{dist}(x, Y))^3} \|t - \eta(t)\|_2 \end{aligned}$$

and this gives the first statement. Applying  $|a + b + c|^2 \leq 3(|a|^2 + |b|^2 + |c|^2)$  gives

$$\begin{aligned} & \frac{\sum_{x_j \notin [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} |(h * (\mu_1 - \mu_2))(x_j)|^2}{\sum_{t \in Y_3} |\tilde{\alpha}_t|^2 + \sum_{t \in Y_1} |\alpha_t - \alpha_{\eta(t)}|^2 + n^2(|\alpha_t|^2 + |\alpha_{\eta(t)}|^2) \|t - \eta(t)\|_2^2} \\ &\leq 3 \max(c_1, \sqrt{2}c_2)^2 n^4 |Y| \sum_{x_j \notin [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} \frac{1}{(1 + n \operatorname{dist}(x_j, Y))^6}. \end{aligned}$$

The remaining sum can be estimated as

$$\sum_{x_j \notin [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} \frac{1}{(1 + n \operatorname{dist}(x_j, Y))^6} \leq \sum_{\ell = \lceil J(\frac{1}{2} + \Delta) \rceil}^{\infty} \sum_{\|j\|_\infty = \ell} \frac{1}{(1 + n(\frac{\ell}{J} - \frac{1}{2}))^6}$$

$$\begin{aligned}
 &= \sum_{\ell=\lceil J(\frac{1}{2}+\Delta) \rceil}^{\infty} \frac{8\ell}{(1+n(\frac{\ell}{J}-\frac{1}{2}))^6} \\
 &= \sum_{\ell=\lceil J(\frac{1}{2}+\Delta) \rceil}^{\infty} \frac{\frac{8J}{n}}{(1+n(\frac{\ell}{J}-\frac{1}{2}))^5} \left(1 + \frac{\frac{n}{2}-1}{1+n(\frac{\ell}{J}-\frac{1}{2})}\right) \\
 &\leq \frac{8J(1+\frac{1}{2\Delta})}{n} \int_{J(\frac{1}{2}+\Delta)}^{\infty} \frac{1}{(1+n(\frac{y}{J}-\frac{1}{2}))^5} dy \\
 &= \frac{8J^2(1+\frac{1}{2\Delta})}{n^2} \int_{\Delta}^{\infty} n(1+nw)^{-5} dw \\
 &= \frac{2J^2(1+\frac{1}{2\Delta})}{n^2} (1+\Delta n)^{-4}
 \end{aligned}$$

and this completes the proof.  $\square$

Taking Lemma 2.2.32 and Lemma 2.2.33 together, we obtain the following result for the condition number of the image recovery problem.

**Theorem 2.2.34.** *Let  $q$  be as in (2.35),  $h$  as in Definition 2.2.29 and  $\Delta > 0$ . There is a constant  $C_{\tau,\gamma,\Delta} > 0$  such that*

$$\frac{\|\tilde{g}^{(1)} - \tilde{g}^{(2)}\|_2^2}{\sum_{j' \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_{j'})|^2} \geq 1 - \frac{C_{\tau,\gamma,\Delta} n^2}{(1+\Delta n)^4}$$

and this lower bound is positive if  $n$  is sufficiently large. From this, we can conclude that there is a constant  $\tilde{c}_{2,\tau,\gamma}^{(5)}$  implying that for  $n$  large enough

$$\tilde{\kappa}_{abs} \left( \frac{\tilde{q}}{n}, \Delta, J, h, M \right) \leq \tilde{c}_{2,\tau,\gamma}^{(5)} \cdot \left( 1 - \frac{C_{\tau,\gamma,\Delta} n^2}{(1+\Delta n)^4} \right)^{-1/2}$$

if  $\tilde{q} = nq = \frac{\sqrt{1+\tau}j_{1,1}}{\pi\gamma} > \frac{j_{1,1}}{\pi}$ .

*Proof.* By (2.34), Lemma 2.2.32 and Lemma 2.2.33, we derive

$$\begin{aligned}
 &\frac{\|\tilde{g}^{(1)} - \tilde{g}^{(2)}\|_2^2}{\sum_{j' \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_{j'})|^2} \\
 &\geq 1 - \frac{\frac{6(1+\frac{1}{2\Delta})J^2|Y|n^2}{(1+\Delta n)^4} \sum_{t \in Y_3} |\tilde{\alpha}_t|^2 + \sum_{t \in Y_1} |\alpha_t - \alpha_{\eta(t)}|^2 + n^2 (|\alpha_t|^2 + |\alpha_{\eta(t)}|^2) \|t - \eta(t)\|_2^2}{4J^2 \max(c_1, \sqrt{2}c_2)^{-2} \left( \min_{\|v\| \leq n'} |\hat{h}(v)|^2 \right) \int_{B_{n'}(0)} |\hat{\mu}_1(v) - \hat{\mu}_2(v)|^2 dv} \\
 &\geq 1 - \frac{3(1+\frac{1}{2\Delta}) \max(c_1, \sqrt{2}c_2)^2 |Y|}{2 \min(c_{2,\tau}^{(1)}, \frac{1}{2}\gamma^2 c_{2,\tau}^{(2)}, c_{2,\tau}^{(3)}) \gamma^2 c_3^2 (1-\gamma)^4 (1+\Delta n)^4} \\
 &\geq 1 - \frac{C_{\tau,\gamma,\Delta} n^2}{(1+\Delta n)^4}
 \end{aligned}$$

where we used  $n' = \gamma n$  for the second and  $|Y| \leq \left( \frac{\sqrt{2}\pi\gamma n}{\sqrt{1+\tau}j_{1,1}} \right)^2$  for the third inequality. For the second part of the theorem, one can at first obtain

$$W_1(\mu_1, \mu_2)^2 \leq \left( \tilde{c}_{d,\tau,\gamma}^{(5)} \right)^2 \sum_{j' \in \mathbb{Z}^2} |(h * (\mu_1 - \mu_2))(x_{j'})|^2$$

## 2 Condition of sparse super resolution

for some  $\tilde{c}_{2,\tau,\gamma}^{(5)} > 0$  by proceeding analogously to the proof of Theorem 2.2.13. Together with the first part of this theorem, one directly concludes the stated upper bound on the condition.  $\square$

After finding this upper bound on the condition of the image recovery problem in the case of sufficiently well-separated sources, it is natural to study again the diffraction limit as the transition point where the condition deteriorates.

**Definition 2.2.35** (Diffraction limit of the imaging problem). We define the *optimal transition constant of the image recovery problem*  $\tilde{\Omega}_2 \geq 0$  as

$$\tilde{\Omega}_2 = \inf \left\{ \tilde{q} > 0 : \exists \beta \in \mathbb{N} \lim_{n \rightarrow \infty} \sup_{M \leq (\sqrt{dn}/\tilde{q})^d} \frac{\tilde{\kappa}_{\text{abs}} \left( \frac{\tilde{q}}{n}, \Delta, J, h, M \right)}{M^\beta} < \infty \right\}.$$

Analogously to the periodic analysis, we obtain the same result on this transition constant as in Theorem 2.2.21. This shows that the switch in the condition is independent of the data being low order Fourier measurements arising in a periodic setting or bandlimited low pass versions of a compactly supported measure.

**Theorem 2.2.36** (Transition constant for the imaging problem). *The transition constant from Definition 2.2.35 satisfies*

$$1.16 \approx \sqrt{\frac{4}{3}} \leq \tilde{\Omega}_2 \leq \frac{j_{1,1}}{\pi} \approx 1.22.$$

*Proof.* The upper bound is a direct consequence of Theorem 2.2.34. For the lower bound, we can take the measures  $\mu_1, \mu_2 \in \mathcal{M}_{+,1}([- \frac{1}{2}, \frac{1}{2}]^2)$  from Lemma 2.2.20 by identifying  $[- \frac{1}{2}, \frac{1}{2}]^2$  with  $\mathbb{T}^2$ . Then, we have for  $q$  satisfying  $nq = \sqrt{\frac{4}{3}}(1 - \epsilon)$

$$\begin{aligned} \tilde{\kappa}_{\text{abs}}(q, \Delta, J, h, M) &= \sup_{\substack{\mu \in \mathcal{M}(q) \\ |Y^\mu| \leq M}} \sup_{\substack{\rho \in \mathbb{C}^{\mathcal{J}} \\ \rho \neq 0}} \inf_{\nu \in \mathcal{R}(\{(h * \mu)(x_j)\}_j + \rho)} \frac{W_1(\nu, \mu)}{\|\rho\|_2} \\ &\geq \frac{W_1(\mu_1, \mu_2)}{\left( \sum_{x_j \in [-\frac{1}{2} - \Delta, \frac{1}{2} + \Delta]^2} |(h * \mu_1)(x_j) - (h * \mu_2)(x_j)|^2 \right)^{1/2}} \\ &\geq \frac{\frac{q}{2}}{4J^2 \int_{B_n(0)} |\mu_1(v) - \mu_2(v)|^2 dv} \\ &\geq \frac{2^{\epsilon(n-1)} \cdot q}{128J^2} \end{aligned}$$

and this grows faster in  $n$  than any polynomial.  $\square$

**Remark 2.2.37** (PSFs). (i) The prototypical example of a bandlimited PSF is the *Airy disc* or *Airy pattern*, e.g. see [2, 16, 75, 28]. It is popularly chosen as it arises naturally by the diffraction of light at a perfectly circular aperture, cf. [2]. Using the bandlimit parameter  $n$  as before, it can be written as

$$h : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad x \mapsto h(x) = n^2 \pi \left( \frac{J_1(n\pi \|x\|_2)}{n\pi \|x\|_2} \right)^2.$$

### 2.3 Application to Vandermonde matrices with pair clusters

Here, the parameter  $n$  is related to the *wavelength*  $\lambda$  of the light and the *numerical aperture* NA by  $n^{-1} = \frac{\lambda}{2\text{NA}}$  such that the quantity  $n^{-1}$  is known as *Abbe's diffraction limit* typically being approximately 200nm for visible light (e.g. cf. [28]). We now want to check the conditions from Definition 2.2.29 for this choice of PSF. As  $h(0) = \frac{1}{2}\pi n^2$  by the series representation of Bessel functions, we can validate assumption (2.31) in Definition 2.2.29 by using the asymptotic expansion from Lemma 1.3.2. Additionally, one can compute the radial derivative of the radial function as

$$\frac{\partial h}{\partial r}(x) = \frac{-2nJ_1(n\pi\|x\|_2)J_2(n\pi\|x\|_2)}{\|x\|_2^2}$$

by using the recurrence relation from Lemma 1.3.2. Together with  $\frac{\partial h}{\partial r}(0) = -\frac{1}{4}\pi^2 n^3$  this gives (2.32). Finally, the Fourier transform of  $h$  which is also called *optical transfer function (OTF)* in microscopy (see [75]) can be computed as

$$\hat{h}(v) = \frac{2}{\pi} \left( \arccos\left(\frac{\|v\|_2}{n}\right) - n^{-1}\|v\|_2\sqrt{1 - n^{-2}\|v\|_2^2} \right)$$

for  $\|v\|_2 \leq n$  and  $\hat{h}(v) = 0$  outside of this ball with radius  $n$ , see [28, Fact 2]. This implies  $\hat{h}(0) = 1$ , i.e. convolution with  $h$  preserves the mass of the measure. Moreover, one can find an integral representation of  $\hat{h}$  by differentiation leading to

$$\hat{h}(v) = \frac{4}{\pi} \int_{\frac{\|v\|_2}{n}}^1 \sqrt{1 - s^2} ds \geq \frac{4}{\pi} \int_{\frac{\|v\|_2}{n}}^1 1 - s ds = \frac{2}{\pi} \left( 1 - \frac{\|v\|_2}{n} \right)^2$$

and this shows (2.33). Therefore, the Airy pattern satisfies the assumptions of Definition 2.2.29.

- (ii) Beyond the theoretical knowledge of bandlimited PSFs like the Airy pattern, experimentally measured PSFs are often acquired in applications where the PSF can not be described in a closed form. For example, an approximation to the PSF is then typically obtained by fitting a multivariate Gaussian into the measurements, e.g. cf. [96, 121, 139]. In the course of the analysis, one discards the high order trigonometric moments where the OTF is small such that the reconstruction of them would be dominated by noise. One could extend the condition analysis of this subsection to point spread functions of this kind if the sampling parameter  $J$  is sufficiently large as the decay conditions from Definition 2.2.29 hold with possibly even stronger decay conditions while the fact, that the PSF is bandlimited, is only used in order to rewrite the sum over  $\mathbb{Z}^2$  as an integral over  $\mathbb{R}^2$  via Lemma 1.2.5. But for this step one could also invoke the regularity of the PSF together with a possibly large sampling parameter  $J$  in order to obtain a lower bound on the sum by the integral. Apart from that, the analysis would not be different to the presented discussion for bandlimited PSFs such that we omit further details at this point.

### 2.3 Application to Vandermonde matrices with pair clusters

As already observed by Diederichs in [38] for  $d = 1$ , the results from the previous subsection yield an estimate for the smallest singular value of the Vandermonde matrix

$$\mathcal{A} = \left( e^{-2\pi i t k} \right)_{k \in \{k \in \mathbb{Z}^d : \|k\|_2 \leq n\}, t \in Y} \quad (2.36)$$

in the situation of a clustered node set  $Y$  where the maximal cluster size is two.

## 2 Condition of sparse super resolution

**Definition 2.3.1** (Cluster and their separation). Let  $Y \subset \mathbb{T}^d$  be a finite set and  $n \in \mathbb{N}$  fixed. We say that two nodes  $t_1, t_2 \in Y$  form a *pair cluster*  $\Lambda = \{t_1, t_2\}$  if  $\|t_1 - t_2\|_{\mathbb{T}^d} < \frac{j_{d/2,1}}{\pi n}$ . Additionally, a set containing only one node is called *single cluster* while we say that  $Y$  admits a *clustered node configuration* if  $Y = \Lambda_1 \cup \Lambda_2 \cup \dots \cup \Lambda_r$  for some  $r \in \mathbb{N}$  and distinct pair or single clusters  $\Lambda_j, j = 1, \dots, r$ . Finally, the *minimal cluster separation* of a clustered node configuration is

$$\text{clusep } Y := \min_{i \neq j} \min_{\substack{t_1 \in \Lambda_i \\ t_2 \in \Lambda_j}} \|t_1 - t_2\|_{\mathbb{T}^d}.$$

We remark that there have been different definitions of a cluster including definitions where nodes are said to form a cluster if they are contained in a box with side length  $n^{-1}$ , e.g. [120]. Our choice to introduce the factor  $\frac{j_{d/2,1}}{\pi}$  is of course due to our proof strategy with the minorant function but it is not artificial to define a cluster with a dimension dependent box side length. A clustered node configuration according to our definition for  $d = 2$  is displayed in Figure 2.10.

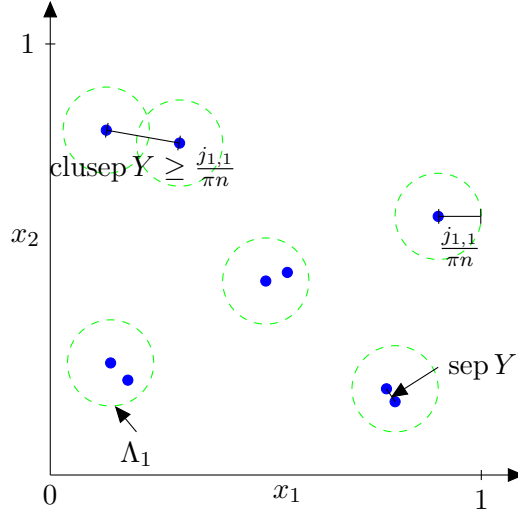


Figure 2.10: Visualisation of pair-cluster configuration in  $d = 2$  where nodes  $t \in Y$  (blue circles) with distance less than  $\frac{j_{1,1}}{\pi n}$  form a cluster and points from different clusters are at least  $\text{clusep } Y \geq \frac{j_{1,1}}{\pi n}$  away from each other. The clusters are highlighted by the green circles.

It is a natural assumption that the clusters need to be well-separated from each other in order to bound the smallest singular value of  $\mathcal{A}$  from below. We make the simple assumption that the clusters are separated by at least the critical separation  $\frac{j_{d/2,1}}{\pi n}$ .

**Proposition 2.3.2** (Pair clusters or well separated nodes). *Assume  $d \in \mathbb{N}$ . Let  $Y$  be a node set with at most pairwise clustering points,  $\text{sep } Y \in (0, \frac{j_{d/2,1}}{\pi n})$  and minimal cluster separation*

$$\text{clusep } Y = \sqrt{1 + \tau} \frac{j_{d/2,1}}{\pi n}$$

for  $\tau > 0$ . Then, a lower bound for the smallest non-zero singular value of the correspond-

### 2.3 Application to Vandermonde matrices with pair clusters

ing Vandermonde matrix  $\mathcal{A}$  in (2.36) is given by

$$\sigma_{\min}(\mathcal{A}) \geq \frac{c_d^{(I)} \sqrt{\tau}}{(1+\tau)^{d/4+1}} \cdot n^{d/2} \cdot \text{SRF}^{-1}$$

for some  $c_d^{(I)} > 0$  and the super resolution factor<sup>49</sup> given by  $\text{SRF} = (n \cdot \text{sep } Y)^{-1}$ . On the contrary, if  $\text{sep } Y = \sqrt{1 + \tau \frac{j_{d/2,1}}{\pi n}}$  for  $\tau \geq 0$ , i.e.  $Y$  is well separated, we have

$$\sigma_{\min}(\mathcal{A}) \geq c_d^{(II)} (1+\tau)^{-d/4-1/2} \cdot n^{d/2}$$

for some dimension dependent constant  $c_d^{(II)} > 0$ .

*Proof.* Let  $\alpha \in \mathbb{C}^{|Y|}$  be an arbitrary normalised vector. We set  $\text{clusep } Y = \sqrt{1 + \tau \frac{j_{d/2,1}}{\pi n}}$  for some  $\tau > 0$  and follow the lines of the proof of Theorem 2.2.8. Next, we decompose the node set into  $Y = Y_4 \cup \bigcup_{t \in Y_5} \{t, \eta(t)\}$  where  $Y_4$  contains all single nodes and  $Y_5$  consists of only one node  $t$  of each pairing node cluster, i.e. the corresponding pair node  $\eta(t)$  is not included in  $Y_5$  in order to circumvent double counting. This allows to compute

$$\begin{aligned} \|\mathcal{A}\alpha\|_2^2 &= \sum_{\|k\|_2 \leq n} \left| \sum_{t \in Y} \alpha_t e^{-2\pi i t k} \right|^2 \\ &\geq \frac{\psi_{\tau,n}(0)}{\hat{\psi}_{\tau,n}(0)} \sum_{t \in Y_4} |\alpha_t|^2 + \sum_{t \in Y_5} \frac{1}{\hat{\psi}_{\tau,n}(0)} \begin{pmatrix} \alpha_t \\ \alpha_{\eta(t)} \end{pmatrix}^* \begin{pmatrix} \psi_{\tau,n}(0) & \psi_{\tau,n}(t - \eta(t)) \\ \psi_{\tau,n}(t - \eta(t)) & \psi_{\tau,n}(0) \end{pmatrix} \begin{pmatrix} \alpha_t \\ \alpha_{\eta(t)} \end{pmatrix} \\ &\geq \min \left( \frac{\psi_{\tau,n}(0)}{\hat{\psi}_{\tau,n}(0)}, \frac{\psi_{\tau,n}(0) - \psi_{\tau,n}(\text{sep } Y \cdot e_1)}{\hat{\psi}_{\tau,n}(0)} \right) \sum_{t \in Y} |\alpha_t|^2. \end{aligned} \quad (2.37)$$

In the clustering case, the second term in the minimum is smaller and the first part follows from the variational formulation of singular values, Lemma 2.2.2 and

$$\frac{\psi_{\tau,n}(0) - \psi_{\tau,n}(\text{sep } Y \cdot e_1)}{\hat{\psi}_{\tau,n}(0)} \geq \frac{c_d \tau (1+\tau)^{-d/2-1} n^2 (\text{sep } Y)^2 n^d}{4\pi^2 (1+\tau) \hat{\varphi}(0)^2} \geq \frac{\left(c_d^{(I)}\right)^2 \tau}{(1+\tau)^{d/2+2}} n^{d+2} (\text{sep } Y)^2$$

for some  $c_d^{(I)} > 0$ . For well-separated nodes, both parts of the minimum in (2.37) are equal and due to the computations in Remark 2.2.9 we obtain

$$\frac{\psi_{\tau,n}(0)}{\hat{\psi}_{\tau,n}(0)} \geq \frac{\left(c_d^{(II)}\right)^2 n^d}{(1+\tau)^{d/2+1}}$$

for some constant  $c_d^{(II)} > 0$ . □

For a result with completely explicit constants for the lower bounds of singular values of Vandermonde matrices which can be derived from Lemma 2.2.6 see [67]. We compare Proposition 2.3.2 to a similar result by Nagel [120, Cor. 3.4.15] which we modified by a frequency shift to fit into our setting.

<sup>49</sup>The term SRF for the super resolution factor is used in many related publications including [11, 100, 108]. The idea is that SRF being smaller or larger than one roughly determines whether one studies weak or strong super resolution. Nevertheless, our analysis from the previous section shows that this distinction should rather be made around  $\frac{\pi}{j_{d/2,1}}$  instead of one.

## 2 Condition of sparse super resolution

**Theorem 2.3.3** (Pair clustering, Nagel’s result). *Let  $Y$  be a node set with at most pairwise clustering points and*

$$\text{clusep } Y \geq \frac{6d}{n} \left( \frac{2}{n \cdot \text{sep } Y} \right)^{\frac{1}{d+1}}.$$

*Then, we have that the smallest non-zero singular value of the corresponding Vandermonde matrix  $\mathcal{A} = (e^{-2\pi i t k})_{k \in \{-n, \dots, n\}^d, t \in Y}$  can be bounded from below by*

$$\sigma_{\min}(\mathcal{A}) \geq \frac{1}{6} \frac{\sqrt{2}}{d^{d/4}} \cdot n^{d/2} \cdot \text{SRF}^{-1}.$$

While contrasting these results for the smallest singular value, we have to remark the following:

- (i) The definition of a cluster is different in the two settings. Whereas in Nagel’s language nodes form a cluster if they are contained in a cube of side length  $\frac{1}{n}$  (independently of the dimension), we considered two nodes as paired if their were separated by less than  $\frac{j_{d/2,1}}{\pi n}$ .
- (ii) The condition on the cluster separation is weaker in Proposition 2.3.2 than in Theorem 2.3.3. Especially, we emphasise that our lower bound for the cluster separation  $\text{clusep } Y$  is independent of  $\text{sep } Y$ , i.e.  $\text{clusep } Y$  does not need to be adjusted for an arbitrarily small  $\text{sep } Y$ .
- (iii) Our result for the smallest singular value has a exponentially better dimension-dependent constant, see [67, Cor. 3.20] where we made the constant  $c_d^{(II)}$  explicit by using the minorant from Lemma 2.2.6.

Consequently, we were able to prove improved lower bounds for the smallest singular values of Vandermonde matrices in the special case of pairwise clustering nodes. Unfortunately, our method does not provide reasonable results for larger clusters containing  $\lambda > 2$  nodes. The latter was done in [120, Sec. 3.4].



### 3 Trigonometric polynomials and rational functions

The first two sections of this chapter are based on the publications [25, 26] while the rest is a previously unpublished extension to approximation by rational functions. Content from [26] about measures with support on algebraic varieties is omitted in this thesis because we focus on super resolution of discrete measures. Only in Section 3.1, the considered measures are allowed to have an arbitrary support.

In this chapter, we present strategies how to use noisy moment information as in (2.2) in order to represent the unknown sparse measure  $\mu$  by a trigonometric polynomial or a rational function which approximate the measure of interest or interpolate its support. Again, the available data consist of trigonometric moments of low to moderate order, i.e.

$$\hat{\mu}(k) = \int_{\mathbb{T}^d} e^{-2\pi i k x} d\mu(x) = \sum_{t \in Y} \alpha_t e^{-2\pi i t \cdot k}, \quad k \in \{-n, \dots, n\}^d \text{ or } k \in \mathbb{Z}^d \cap B_n(0) \quad (3.1)$$

for some  $n \in \mathbb{N}$ , and one asks for the reconstruction or approximation of  $\mu \in \mathcal{M}(q)$  from this partial information. For the approximation using convolutions with kernels in Section 3.1, it is easier to assume knowledge of moments for  $k \in \{-n, \dots, n\}^d$  as we did in [26], while we deal with the *radial setting*  $k \in \mathbb{Z}^d \cap B_n(0)$  in the rest of this chapter. This radial setting is more reasonable in applications as we have explained in Chapter 2 and allows to use the condition estimates from this previous chapter. As the full recovery of all parameters of  $\mu$  is often beyond the scope of the application due to the large number of parameters or because one just wants a good visualisation of the measure, we propose trigonometric proxies  $q_n$  based on the knowledge of (3.1) and distinguish between pointwise convergence to the indicator function of  $\text{supp } \mu$ , i.e.

$$\lim_{n \rightarrow \infty} q_n(x) = \begin{cases} 1, & x \in \text{supp } \mu, \\ 0, & \text{else,} \end{cases} \quad (3.2)$$

and weak convergence  $q_n \rightharpoonup \mu$  as introduced in Section 1.4. We quantify convergence rates for the latter by the Wasserstein distance of  $\mu$  and  $q_n$ .

**Related work and contributions** We have already described in the context of the analysis of the condition of (3.1) that there is a wide range of algorithms which compute the parameters of  $\mu \in \mathcal{M}(q)$  given the trigonometric moments. Beginning with Prony’s method [79, 89, 133, 136, 142] one could mention other subspace methods such as matrix pencil [47, 74, 117], ESPRIT [4, 99, 137, 140] or MUSIC [101, 144]. Among them, MUSIC as well as the variational methods [22, 31, 32, 134] set up intermediate trigonometric polynomials which peak around the support points and have smaller value apart from the support. Beyond the consideration of polynomials, rational functions with this property are studied as well, for instance Christoffel functions offer interesting guarantees both in terms of support identification [93] or approximation *on the support* [84, 111, 128]. Hence, it is

### 3 Trigonometric polynomials and rational functions

natural to ask whether these different polynomial or rational proxies interpolate the support of a discrete measure and fulfil a pointwise convergence similar to (3.2). In contrast to this approximation of the support by interpolation, weak convergence of polynomial approximations to the measure itself is studied in particular in Mhaskar’s paper [113]. Given the moments up to order  $n$ , we follow a similar approach by addressing two types of questions:

- (QA) Which simply computable, trigonometric approximation schemes for general measures are available and how well do they approximate the measure?
- (QB) How close are these schemes to the actual best polynomial approximation?

In contrast to [113], the discrepancy between the measure and its polynomial proxy is estimated in the 1-Wasserstein distance which gives rise to tight bounds on the approximation error. One of the main contributions of our work lies in the simple connection between approximation in the 1-Wasserstein distance and known results from approximation theory for Lipschitz functions. For example, we relate questions on the best approximation of measures by polynomials to best approximation results in  $L^1(\mathbb{T}^d)$  and  $C(\mathbb{T}^d)$ . Additionally, we show analogously to classical approximation theory that near best approximations can be derived through convolution with certain kernels. As far as we know, these connections formulated in Section 3.1 were not considered before.

For the pointwise convergence (3.2), we analyse in Section 3.2 the interpolation behaviour of a sum of squares polynomial,  $p_{1,n}$ , similarly suggested in [89, Thm. 3.5] and [124, Prop. 5.3] (and indeed closely related to the *rational* function in the MUSIC algorithm, see [144, Eq. (6)]). The main contribution of this section is not the invention of this polynomial  $p_{1,n}$  but the analysis of its pointwise convergence to the indicator function of the support of the measure. For instance, this can be used in Section 4.1 to represent sparse objects in single molecule microscopy.

**Organisation of this chapter** Section 3.1 presents our answers towards (QA) and (QB). Similar to classical approximation theory for functions, the simplest approximation methods are given by convolution of the measure with polynomial kernels and we derive upper bounds in the 1-Wasserstein distance in Subsection 3.1.1. Comparing these results for question (QA) with the analysis of best approximations in Subsection 3.1.3, we observe for question (QB) that the Fejér kernel provides a sub-optimal order whereas the Jackson kernel achieves the optimal rate of  $\mathcal{O}(n^{-1})$ . Furthermore, the subsequent characterisation of the best approximation in the univariate case in Subsection 3.1.4 shows that the achieved constant in the convergence rate for the best approximation is sharp.

In Section 3.2, we start by studying the so-called signal polynomial  $p_{1,n}$  which identifies the support of a discrete measure in the sense of (3.2). As common to all subspace methods, this involves technical assumptions on the support of the measure and the degree  $n$  to be finite but large enough. For discrete measures with separation larger than the Rayleigh limit, see Subsection 2.2.3, we can apply results from Chapter 2 and prove in Subsection 3.2.2, Theorem 3.2.6, the pointwise convergence to the indicator function of support  $\mu$  with a pointwise convergence rate of order  $\mathcal{O}(n^{-2})$  outside of the support.

Beyond studying representations by polynomials, we propose to use rational approximations in Section 3.3. Naturally, the situation becomes more involved if we include noise into our data model (3.1) and an approach how to adjust the rational approximation approach in this setting is presented in Section 3.4.

## 3.1 Approximation by polynomials

We study in this section weakly convergent polynomial approximations of measures, i.e. approximations satisfying the property (1.9). While general  $p$ -Wasserstein distances allow to quantify this weak convergence with a convergence rate, see [151], we restrict ourselves to the 1-Wasserstein distance. Some bounds in the setting of  $W_p$ ,  $p > 1$ , were derived in [25] but are not included in this thesis completely in order to avoid redundancies.<sup>50</sup>

While our focus is principally on actually computable approximations, based on convolution with known kernels, we also turn in the last part of this section (Subsection 3.1.3 below) to more theoretical considerations on the best polynomial approximations with respect to the 1-Wasserstein distance, which, additionally to giving new perspectives on polynomial approximations of measures, also highlights the near-optimality of our constructions.

### 3.1.1 Approximation by convolution and upper bounds

Similarly to standard approaches in approximation theory, one may derive easy-to-compute polynomial estimates for a measure  $\mu$ , by considering the convolution of the latter with adequate kernels. As highlighted in the beginning of this chapter, we simplify the setting in contrast to Chapter 2 by considering tensor product kernels for trigonometric moments from the box  $\{-n, \dots, n\}^d$  instead of a radial problem with moments corresponding to frequencies from the Euclidean ball  $B_n(0)$  because the latter would enforce to study kernels from  $\mathcal{P}^{n,d,2}$  which is less straightforward even though the results can only differ by a dimension dependent constant due to  $\mathcal{P}^{d-1/2,n,d,\infty} \subset \mathcal{P}^{n,d,2} \subset \mathcal{P}^{n,d,\infty}$ .

Given the first trigonometric moments  $(\hat{\mu}(k))_{k \in \{-n, \dots, n\}^d}$  of  $\mu$ , a first natural choice of a kernel would be to take the Fourier partial sums

$$S_n \mu(x) = (D_n * \mu)(x) = \sum_{k \in \mathbb{Z}^d, \|k\|_\infty \leq n} \hat{\mu}(k) e^{2\pi i k x},$$

which correspond to convolution with Dirichlet kernels, cf. Section 1.3. We focus in this section on yet another classical sequence of approximations, given by convolution with Fejér kernels  $F_n: \mathbb{T}^d \rightarrow \mathbb{R}$  (by slight abuse of notation, we use the same notation for both the multivariate and univariate kernels). We recall that its definition, see Definition 1.3.5, allows to write

$$F_n(x) = \sum_{k=-n}^n \left(1 - \frac{|k|}{n+1}\right) e^{2\pi i k x} = \frac{1}{n+1} \left(\frac{\sin(n+1)\pi x}{\sin \pi x}\right)^2 = \frac{1}{n+1} \left|\sum_{k=0}^n e^{2\pi i k x}\right|^2 \quad (3.3)$$

for any  $x \in \mathbb{T}$ . The main object of study in this section is the trigonometric polynomial

$$p_n(x) = (F_n * \mu)(x) = \int_{\mathbb{T}^d} F_n(x-y) d\mu(y). \quad (3.4)$$

We start by giving two illustrative examples in Example 3.1.1.

**Example 3.1.1.** Our first example for  $d = 1$  is the measure

$$\mu = \frac{1}{3} \delta_{\frac{1}{8}} + \nu \in \mathcal{M}_+(\mathbb{T}), \quad \frac{d\nu}{d\lambda}(x) = \frac{8}{9} \mathbb{1}_{[\frac{1}{4}, \frac{5}{8}]}(x) + \frac{\sqrt{2}}{3} \left( \frac{1}{\sqrt{|x - \frac{7}{8}|}} - \sqrt{8} \right) \mathbb{1}_{[\frac{3}{4}, 1]}(x), \quad (3.5)$$

<sup>50</sup>At some points we will remark about the behaviour of certain bounds for general  $p$  as presented in [25].

### 3 Trigonometric polynomials and rational functions

where  $\lambda$  denotes the Lebesgue measure. Obviously, the measure  $\mu$  has singular and absolutely continuous parts including an integrable pole at  $x = \frac{7}{8}$ . Both the Fourier partial sums and the Fejér approximations for  $n = 19$  are shown in the left and right panel of Figure 3.1, respectively.

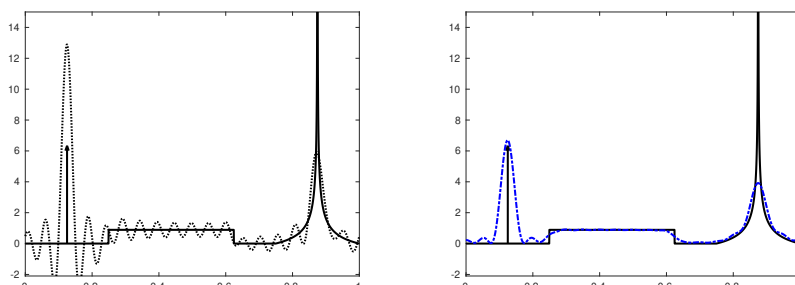


Figure 3.1: The example measure (3.5), its approximations by the Fourier partial sum (left) and the convolution with the Fejér kernel (right). The weight  $\frac{1}{3}$  of the Dirac measure in  $\frac{1}{8}$  is displayed by an arrow of height  $n/3$  for visibility.

Our second example is a singular continuous measure for  $d = 2$ . We take the measure  $\mu = (2\pi r_0)^{-1} \delta_C \in \mathcal{M}_+(\mathbb{T}^2)$  as the uniform measure on the circle

$$C = \{x \in \mathbb{T}^2 : |x|_2 = r_0\}$$

for some radius  $0 < r_0 < \frac{1}{2}$ . The total variation of this measure is

$$\|\mu\|_{\text{TV}} = \hat{\mu}(0) = \int_{\mathbb{T}^2} d\mu(x) = \frac{1}{2\pi r_0} \int_C dx = 1.$$

As we have  $\text{supp } \mu \subset [-\frac{1}{2}, \frac{1}{2}]^2$ , we find with (1.1) the explicit formula

$$\hat{\mu}(k) = \int_{\mathbb{T}^2} e^{-2\pi i k x} d\mu(x) = \frac{1}{r_0} \int_0^\infty r J_0(2\pi r \|k\|_2) d\delta_{r_0}(r) = J_0(2\pi r_0 \|k\|_2) \quad (3.6)$$

for the trigonometric moments of  $\mu$ . The Fourier partial sum as well as the convolution with the Fejér kernel for  $n = 29$  are shown with maximal contrast in the left and right panel

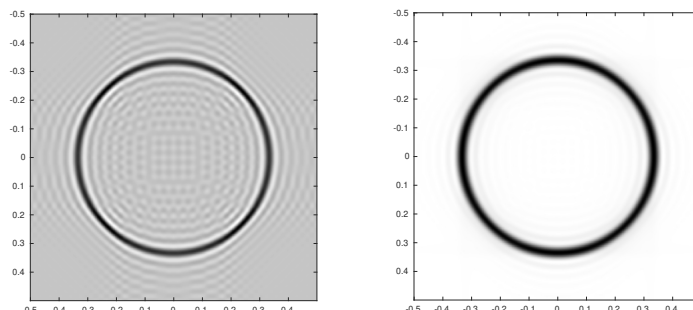


Figure 3.2: Uniform measure on a circle of radius  $r_0 = \frac{1}{3}$  and its approximations by the Fourier partial sum (left) and the convolution with the Fejér kernel (right) for  $n = 29$ .

of Figure 3.2, respectively. We observe that both approximators peak around the support of the measures and the approximation by convolution with the Fejér kernel produces less ringing than the Dirichlet kernel at the cost of a slightly thicker main lobe. Especially the convolution with the Fejér kernel can be seen as a visually appealing proxy for  $\mu$ .

Of course, the construction and efficient evaluation of this approximation relies on the convolution theorem and the fast Fourier transform (FFT). Given the trigonometric moments  $\hat{\mu}(k)$ ,  $k \in \{-n, \dots, n\}^d$ , we multiply these with the Fourier coefficients of the Fejér kernel (3.3) in each dimension and use an inverse FFT to evaluate  $p_n$  on the equispaced grid  $(2n+1)^{-1}\{-n, \dots, n\}^d$ . Our next goal is a quantitative approximation result, for which we need the following preparatory lemma. This result can be found in qualitative form e.g. in [20, Lemma 1.6.4]. Here, we let  $|x-y|_p = \min_{k \in \mathbb{Z}^d} \|x-y+k\|_p$  for  $p \geq 1$  denote the wrap-around  $p$ -norm on  $\mathbb{T}^d = [0, 1)^d$  such that we can rewrite the metric on  $\mathbb{T}^d$  as  $\|x-y\|_{\mathbb{T}^d} = |x-y|_2$ .

**Lemma 3.1.2.** *Let  $n, d \in \mathbb{N}_{\geq 1}$ , then we have*

$$\frac{d}{\pi^2} \left( \frac{\log(n+1)}{n+1} + \frac{1}{n+2} \right) \leq \int_{\mathbb{T}^d} F_n(x) |x|_1 dx \leq \frac{d}{\pi^2} \frac{\log(n)+4}{n+1}.$$

*Proof.* First note that

$$\int_{\mathbb{T}^d} \prod_{s=1}^d F_n(x_s) \sum_{\ell=1}^d |x_\ell|_1 dx = \sum_{\ell=1}^d \int_{\mathbb{T}^d} \prod_{s=1}^d F_n(x_s) |x_\ell|_1 dx = d \int_{\mathbb{T}} F_n(x) |x|_1 dx,$$

where the second equality holds since  $\int F_n(x_s) dx_s = 1$ . Thus it is sufficient to consider the univariate case. The representation  $F_n(x) = 1 + 2 \sum_{k=1}^n \left(1 - \frac{k}{n+1}\right) \cos(2\pi kx)$  gives

$$\begin{aligned} \int_{\mathbb{T}} F_n(x) |x|_1 dx &= 2 \int_0^{1/2} \left( x + 2 \sum_{k=1}^n \left(1 - \frac{k}{n+1}\right) \cos(2\pi kx) x \right) dx \\ &= 2 \left[ \frac{1}{8} + \sum_{k=1}^n \frac{(-1)^k - 1}{2\pi^2 k^2} - \sum_{k=1}^n \frac{(-1)^k - 1}{2(n+1)\pi^2 k} \right] \\ &= 2 \left[ \frac{1}{8} - \sum_{j=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{1}{\pi^2 (2j+1)^2} + \sum_{j=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{1}{(n+1)\pi^2 (2j+1)} \right] \end{aligned}$$

since  $\int_0^{1/2} \cos(2\pi kx) x dx = ((-1)^k - 1)/(4\pi^2 k^2)$ . Using that  $\sum_{j=0}^{\infty} \frac{1}{(2j+1)^2} = \frac{\pi^2}{8}$ , we obtain

$$\begin{aligned} \int_{\mathbb{T}} F_n(x) |x|_1 dx &= 2 \left[ \frac{1}{\pi^2} \sum_{j=\lfloor \frac{n+1}{2} \rfloor}^{\infty} \frac{1}{(2j+1)^2} + \frac{1}{(n+1)\pi^2} \sum_{j=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{1}{2j+1} \right] \\ &\leq \frac{2}{\pi^2} \left[ \frac{1}{(2\lfloor \frac{n+1}{2} \rfloor + 1)^2} + \int_{\lfloor \frac{n+1}{2} \rfloor}^{\infty} \frac{1}{(2y+1)^2} dy + \frac{1 + \int_0^{\lfloor \frac{n-1}{2} \rfloor} \frac{1}{2y+1} dy}{n+1} \right] \\ &\leq \frac{\frac{2}{(2\lfloor \frac{n+1}{2} \rfloor + 1)^2} + 1}{(2\lfloor \frac{n+1}{2} \rfloor + 1)\pi^2} + \frac{2 + \log(n)}{(n+1)\pi^2} \leq \frac{\log(n)+4}{\pi^2(n+1)}. \end{aligned}$$

### 3 Trigonometric polynomials and rational functions

The lower bound follows similarly by bounding the series from the previous calculation by integrals from below.  $\square$

**Theorem 3.1.3.** *Let  $d, n \in \mathbb{N}$  and  $\mu \in \mathcal{M}(\mathbb{T}^d)$ , then the measure with density  $p_n$  from (3.4) converges weakly to  $\mu$  with*

$$W_1(p_n, \mu) \leq \frac{d \log(n) + 4}{\pi^2 (n + 1)} \cdot \|\mu\|_{\text{TV}},$$

which is almost sharp<sup>51</sup> since we have

$$\sup_{\mu \in \mathcal{M}(\mathbb{T}^d) \setminus \{0\}} \frac{W_1(p_n, \mu)}{\|\mu\|_{\text{TV}}} \geq d^{-1/2} \frac{d}{\pi^2} \left( \frac{\log(n+1)}{n+1} + \frac{1}{n+2} \right).$$

*Proof.* We compute

$$\begin{aligned} W_1(p_n, \mu) &= \sup_{\text{Lip}(f) \leq 1} |\langle F_n * \mu, f \rangle - \langle \mu, f \rangle| \\ &= \sup_{\text{Lip}(f) \leq 1} |\langle \mu, F_n * f - f \rangle| \\ &\leq \sup_{\text{Lip}(f) \leq 1} \int_{\mathbb{T}^d} \int_{\mathbb{T}^d} F_n(x) |f(y-x) - f(y)| dx d\mu(y) \\ &\leq \|\mu\|_{\text{TV}} \int_{\mathbb{T}^d} F_n(x) |x|_2 dx, \end{aligned}$$

and note that both inequalities become equalities when choosing  $\mu = \delta_0$  and  $f(x) = |x|_2$ . Applying  $|\cdot|_2 \leq |\cdot|_1$  and Lemma 3.1.2 gives the first part of the result. In particular, we remark in passing that  $W_1(F_n, \delta_0) = \int_{\mathbb{T}^d} F_n(x) |x|_2 dx \geq d^{-1/2} \int_{\mathbb{T}^d} F_n(x) |x|_1 dx$  yields the second part.  $\square$

**Remark 3.1.4.** In [25], we have analysed that the above weak convergence is only of order  $\mathcal{O}(n^{-1/p})$  if we measure in  $W_p$  for  $p > 1$ . Similar to classical results from approximation theory, the log-factor in Theorem 3.1.3 can be removed by choosing another convolution kernel, which then however does not allow for the representation later found in Lemma 3.2.1. For example, the Jackson kernel  $J_n$  as introduced in Definition 1.3.6 has degree  $n = 2m - 2$  and satisfies

$$\int_{\mathbb{T}} J_n(x) |x|_1 dx \leq \frac{6}{m(2m^2 + 1)} \left[ \int_0^{1/2m} m^4 x dx + \int_{1/2m}^\infty \frac{1}{16x^3} dx \right] \leq \frac{3m}{4m^2 + 2} \leq \frac{3}{2(n+2)}.$$

Analogously to Theorem 3.1.3, we get

$$W_1(J_n * \mu, \mu) \leq \frac{3d \cdot \|\mu\|_{\text{TV}}}{2(n+2)}, \quad (3.7)$$

which still is an approximate factor 6 worse than the lower bound in the univariate case (see Theorem 3.1.6 in Subsection 3.1.3). By numerical analysis or more detailed analysis of the above estimate, one can deduce that a factor 3 is due to the above estimate and a remaining factor 2 seems to indicate that the Jackson kernel is not optimal.

<sup>51</sup>In [26], we defined the Wasserstein distance with respect to  $|\cdot|_1$  on  $\mathbb{T}^d$  such that we can get rid of the factor  $d^{-1/2}$  in the lower bound and obtain an actually sharp bound. In order to be consistent with the rest of this thesis, we stick to our original definition of  $W_1$  which gives this almost sharp lower bound.

### 3.1.2 Saturation

Theorem 3.1.3 gives a worst case lower bound while, on the other hand, the Lebesgue measure is approximated by  $F_n * \lambda = \lambda$  without any error. We may thus ask how well a measure  $d\mu = w(x)dx$  with smooth (nonnegative) density might be approximated. For an introductory example, consider the univariate analytical density  $w(x) = 1 + \cos(2\pi x)$ . Since  $F_n * w(x) - w(x) = \cos(2\pi x)/(n+1)$ , by testing with the Lipschitz function  $f(x) = \cos(2\pi x)/(2\pi)$ , we achieve

$$W_1(F_n * w, w) \geq \frac{1}{2\pi(n+1)} \int_{\mathbb{T}} \cos^2(2\pi x) dx = \frac{1}{4\pi(n+1)}.$$

This effect is called saturation (e.g. cf. [20]). In greater generality, such a lower bound holds for each measure individually and can be inferred by a nice relationship between the Wasserstein distance and a discrepancy, cf. [46].

**Theorem 3.1.5.** *For each individual measure  $\mu \in \mathcal{M}(\mathbb{T}^d)$  different from the Lebesgue measure, there is a constant  $c > 0$  such that*

$$W_1(p_n, \mu) \geq \frac{c}{n+1}$$

holds for all  $n \in \mathbb{N}$ .

*Proof.* Let  $\hat{h} \in \ell^2(\mathbb{Z}^d)$ ,  $\hat{h}(k) \in \mathbb{R} \setminus \{0\}$ ,  $\hat{h}(k) = \hat{h}(-k)$ , and consider the reproducing kernel Hilbert space

$$H = \left\{ f \in L^2(\mathbb{T}^d) : \sum_{k \in \mathbb{Z}^d} |\hat{h}(k)|^{-2} |\hat{f}(k)|^2 < \infty \right\}, \quad \|f\|_H^2 = \sum_{k \in \mathbb{Z}^d} |\hat{h}(k)|^{-2} |\hat{f}(k)|^2.$$

Given measures  $\mu, \nu$ , their discrepancy (which also depends on the space  $H$ ) is defined by

$$\mathcal{D}(\mu, \nu) = \sup_{\|f\|_H \leq 1} \left| \int_{\mathbb{T}^d} f d(\mu - \nu) \right| = \sup_{\|f\|_H \leq 1} \left| \sum_{k \in \mathbb{Z}^d} \frac{\hat{f}(k)}{\hat{h}(k)} \widehat{h(k)\mu - \nu}(k) \right| = \|\widehat{h} \cdot \widehat{\mu - \nu}\|_{\ell^2},$$

and fulfils by the geometric-arithmetic inequality

$$\begin{aligned} \mathcal{D}(p_n, \mu)^2 &= \sum_{\|k\|_\infty \leq n} |\hat{h}(k)|^2 \left| 1 - \prod_{\ell=1}^d \left( 1 - \frac{|k_\ell|}{n+1} \right) \right|^2 |\hat{\mu}(k)|^2 + \sum_{\|k\|_\infty > n} |\hat{h}(k)|^2 |\hat{\mu}(k)|^2 \\ &\geq \sum_{\|k\|_\infty \leq n} |\hat{h}(k)|^2 \left| \frac{\|k\|_1}{d(n+1)} \right|^2 |\hat{\mu}(k)|^2 + \sum_{\|k\|_\infty > n} |\hat{h}(k)|^2 |\hat{\mu}(k)|^2 \\ &= \sum_{\|k\|_\infty \leq n} |\hat{h}(k)|^2 \left| \frac{\|k\|_1}{d(n+1)} \right|^2 |\hat{\mu}(k) - \hat{\lambda}(k)|^2 + \sum_{\|k\|_\infty > n} |\hat{h}(k)|^2 |\hat{\mu}(k) - \hat{\lambda}(k)|^2 \\ &\geq \frac{1}{d^2(n+1)^2} \|h * (\mu - \lambda)\|_{L^2(\mathbb{T}^d)}^2 \end{aligned}$$

where  $h(x) = \sum_{k \in \mathbb{Z}^d} \hat{h}(k) e^{2\pi i k x}$  and  $\lambda$  denotes the Lebesgue measure with  $\hat{\lambda}(0) = 1$  and  $\hat{\lambda}(k) = 0$  for  $k \in \mathbb{Z}^d \setminus \{0\}$ . Our second ingredient is a Lipschitz estimate that quantifies

### 3 Trigonometric polynomials and rational functions

the Lipschitz constant of any  $f \in H$  with  $\|f\|_H \leq 1$ . For such a function  $f$ , the Cauchy-Schwarz inequality together with  $|e^{2\pi iky} - e^{2\pi ik(y+x)}|^2 = 2(1 - \cos(2\pi kx))$  gives

$$\begin{aligned} |f(y) - f(y+x)|^2 &= \left| \sum_{k \in \mathbb{Z}^d} \hat{f}(k) \left( e^{2\pi iky} - e^{2\pi ik(y+x)} \right) \right|^2 \\ &\leq \|f\|_H^2 \sum_{k \in \mathbb{Z}^d} |e^{2\pi iky} - e^{2\pi ik(y+x)}|^2 |\hat{h}(k)|^2 \\ &= \|f\|_H^2 \sum_{k \in \mathbb{Z}^d} (2 - 2\cos(2\pi kx)) |\hat{h}(k)|^2 \\ &\leq 2(K(x, x) - K(x, 0)), \end{aligned}$$

where  $K(x, y) = \sum_{k \in \mathbb{Z}^d} |\hat{h}(k)|^2 e^{2\pi ik(x-y)} = (h * h)(x-y)$  denotes the so-called reproducing kernel of the space  $H$ .<sup>52</sup> If this kernel is  $K(x, y) = h^{[4]}(x_1 - y_1) \cdot \dots \cdot h^{[4]}(x_d - y_d)$  for some nonnegative univariate function  $h^{[4]} \in C^2(\mathbb{T})$  being maximal in zero (and thus  $(h^{[4]})'(0) = 0$ ), we find by a telescoping sum and the Taylor expansion

$$\begin{aligned} K(x, x) - K(x, 0) &= \prod_{\ell=1}^d h^{[4]}(0) - \prod_{\ell=1}^d h^{[4]}(x_\ell) \\ &= \sum_{\ell=1}^d \left( h^{[4]}(0)^\ell \prod_{k=1}^{d-\ell} h^{[4]}(x_k) - h^{[4]}(0)^{\ell-1} \prod_{k=1}^{d-\ell+1} h^{[4]}(x_k) \right) \\ &= \sum_{\ell=1}^d \left( h^{[4]}(0)^{\ell-1} \left( h^{[4]}(0) - h^{[4]}(x_{d-\ell+1}) \right) \prod_{k=1}^{d-\ell} h^{[4]}(x_k) \right) \\ &\leq \sum_{\ell=1}^d \|h^{[4]}\|_\infty^{d-1} \left[ h^{[4]}(0) - h^{[4]}(x_{d-\ell+1}) \right] \\ &\leq \frac{1}{2} \|h^{[4]}\|_\infty^{d-1} \left\| \left( h^{[4]} \right)'' \right\|_\infty |x|_2^2. \end{aligned}$$

To make a specific choice, let  $a \in (0, \frac{1}{8})$  be some irrational number and set

$$h^{[2]}(x) = \sum_{k \in \mathbb{Z}} (\mathbb{1}_{[-a, a]} * \mathbb{1}_{[-a, a]})(x+k), \quad x \in \mathbb{T}$$

as the periodisation of the convolution of the indicator function on  $[-a, a]$  with itself. Based on this, we set  $h^{[4]} = h^{[2]} * h^{[2]}$ , and  $h(x_1, \dots, x_d) = h^{[2]}(x_1) \cdot \dots \cdot h^{[2]}(x_d)$ .<sup>53</sup> Consequently,

<sup>52</sup>Note that the assumptions  $\hat{h}(k) \in \mathbb{R} \setminus \{0\}$  and  $\hat{h}(k) = \hat{h}(-k)$  lead to  $K(x, y) = \sum_{k \in \mathbb{Z}^d} |\hat{h}(k)|^2 e^{2\pi ik(x-y)} = \sum_{k \in \mathbb{Z}^d} |\hat{h}(k)|^2 \cos(2\pi k(x-y))$  and in particular  $K$  is real valued.

<sup>53</sup>Note that the Fourier coefficients of  $h^{[2]}$  agree with the Fourier transform of  $\mathbb{1}_{[-a, a]} * \mathbb{1}_{[-a, a]}$  evaluated at integers by the Poisson summation formula, and analogously this holds for  $h^{[4]}$  and the higher order spline obtained by threefold convolution of  $\mathbb{1}_{[-a, a]}$  with itself. By choosing  $a < \frac{1}{8}$ ,  $h^{[2]}$  and  $h^{[4]}$  agree with these compactly supported convolutions on  $[-\frac{1}{2}, \frac{1}{2}]$ . One immediately gets  $\widehat{h^{[4]}}(k) \in \mathcal{O}(k^{-4})$  by the convolution theorem of the Fourier transform and this indeed yields  $h^{[4]} \in C^2(\mathbb{T})$  by [61, Prop. 3.3.12 or Ex. 2.4.1] meaning that the choice of  $h^{[4]}$  is compatible with our previous assumptions on it. Moreover, we directly have summability of  $|\hat{h}(k)|^2$  for  $k \in \mathbb{Z}^d$  such that  $K$  is a valid kernel.



we derive that  $f \in H$  with  $\|f\|_H \leq 1$  satisfies  $\text{Lip}(f) \leq c'_{d,a}$  for some constant  $c'_{d,a} > 0$  depending on the dimension  $d$  and the parameter  $a$ . This allows to conclude that

$$W_1(p_n, \mu) \geq c'^{-1}_{d,a} \mathcal{D}(p_n, \mu) \geq \frac{c'^{-1}_{d,a}}{d(n+1)} \|h * (\mu - \lambda)\|_{L^2(\mathbb{T}^d)} =: \frac{c}{n+1}$$

for some  $c \in \mathbb{R}$ . Since  $a$  is irrational, we can see by Parseval's theorem that  $\|h * (\mu - \lambda)\|_{L^2(\mathbb{T}^d)} = 0$  if and only if  $\mu = \lambda$ .<sup>54</sup> For  $\mu \neq \lambda$ , we obtain the statement with a positive constant  $c$  depending on the measure  $\mu$ , the constant  $a$ , and the spatial dimension  $d$ .  $\square$

### 3.1.3 Best approximation and lower bounds

After observing upper (Subsection 3.1.1) and lower bounds on the approximation by  $p_n = F_n * \mu$  for individual measures  $\mu$  (Subsection 3.1.2), one might ask whether an approximation rate faster than  $\mathcal{O}(n^{-1})$  is possible by some general polynomial approximation, see (QB) from the beginning of this chapter. The following theorem shows that the answer to this question is negative as the best approximation by a normalised polynomial only yields a  $\mathcal{O}(n^{-1})$  worst-case rate. In order to simplify the problem, we will assume in this and the following subsection that  $\mu$  and its polynomial proxy  $p$  with  $p(x) = \sum_{k \in \{-n, \dots, n\}^d} \mathbf{p}_k e^{2\pi i k x}$  are normalised, i.e.  $\hat{\mu}(0) = \mathbf{p}_0 = 1$ .

**Theorem 3.1.6.** *For any  $d, n \in \mathbb{N}$ ,  $n \geq 1$  and for every  $\mu \in \mathcal{M}(\mathbb{T}^d)$  there exists a polynomial with degree  $n$  of best approximation in the 1-Wasserstein distance among all polynomials with degree  $n$ . Moreover, we have*

$$\sup_{\substack{\mu \in \mathcal{M}(\mathbb{T}^d) \\ \mu(\mathbb{T}^d) = 1}} \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} \frac{W_1(p, \mu)}{\|\mu\|_{\text{TV}}} \geq \frac{d^{-1/2}}{4(n+1)}.$$

*Proof.* We directly have existence of a best approximation by polynomials in the Banach space of Borel measures with finite total variation (e.g. cf. [36, Thm. 3.1.1]). For the lower bound, we compute

$$\begin{aligned} \sup_{\substack{\mu \in \mathcal{M}(\mathbb{T}^d) \\ \mu(\mathbb{T}^d) = 1}} \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} \frac{W_1(p, \mu)}{\|\mu\|_{\text{TV}}} &\geq \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} W_1(p, \delta_0) \\ &= \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} \sup_{\substack{f: \|f\|_\infty \leq \frac{\sqrt{d}}{4} \\ \text{Lip}(f) \leq 1}} \left| f(0) - \int_{\mathbb{T}^d} f(x) p(x) dx \right| \\ &= \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} \sup_{\substack{f: \|f\|_\infty \leq \frac{\sqrt{d}}{4} \\ \text{Lip}(f) \leq 1}} \|f - \check{p} * f\|_\infty \\ &\geq \sup_{\substack{f: \|f\|_\infty \leq \frac{\sqrt{d}}{4} \\ \text{Lip}(f) \leq 1}} \min_{p \in \mathcal{P}^{n,d,\infty}} \|f - p\|_\infty, \end{aligned}$$

where  $\check{p}$  denotes the *reflection* of  $p$ , i.e.  $\check{p}(x) = p(-x)$  for all  $x \in \mathbb{T}^d$ . It remains to find the worst case error for the best approximation of a Lipschitz function by a trigonometric polynomial of degree  $n$ . While this is well-understood for  $d = 1$  (cf. [3, 51]), we did not

<sup>54</sup>We remark that  $\hat{h}(k) = \prod_{\ell=1}^d \frac{\sin^2(\pi k_\ell a)}{\pi^2 k_\ell^2} \neq 0$  for  $a$  irrational. Hence,  $\|h * (\mu - \lambda)\|_{L^2(\mathbb{T}^d)} = 0$  implies by Parseval's theorem that  $\hat{\mu}(k) = \hat{\lambda}(k)$  for any  $k \in \mathbb{Z}^d$ . The latter is equivalent to  $\mu = \lambda$ .

### 3 Trigonometric polynomials and rational functions

find a reference mentioning whether and how  $d > 1$  is possible as well. Therefore, we show that the idea by [55] for the case  $d = 1$  works also for  $d > 1$  in our situation. A main ingredient of Fishers proof is the duality relation

$$\inf_{x \in Y \subset X} \|x_0 - x\| = \sup_{\substack{\ell \in X^* \\ \ell|_Y=0, \|\ell\|_{X^*} \leq 1}} |\ell(x_0)|$$

for a Banach space  $X$ ,  $x_0 \in X$ , with a subspace  $Y$  and dual space  $X^*$ . A second ingredient is given by the 1-periodic Bernoulli spline of degree 1

$$\mathcal{B}_1(x) = \sum_{k \in \mathbb{Z} \setminus \{0\}} \frac{e^{2\pi i k x}}{2\pi i k} = \sum_{k=1}^{\infty} \frac{\sin(2\pi k x)}{\pi k} = \begin{cases} \frac{1}{2} - x, & x \in (0, 1), \\ 0, & x \in \{0, 1\} \end{cases} \quad (3.8)$$

for  $x \in [0, 1]$ .<sup>55</sup> A Lipschitz continuous and 1-periodic function  $f: \mathbb{T} \rightarrow \mathbb{R}$  with  $\text{Lip}(f) \leq 1$  has a derivative  $f'$  almost everywhere<sup>56</sup> and this derivative satisfies  $\int_{\mathbb{T}} f'(s) ds = 0$  by the periodicity of  $f$ . Therefore, it follows that

$$\begin{aligned} (f' * \mathcal{B}_1)(t) &= \int_{\mathbb{T}} f'(s) \mathcal{B}_1(t-s) ds \\ &= - \int_0^t (t-s) f'(s) ds - \int_t^1 (t-s+1) f'(s) ds \\ &= f(t) - \int_0^1 f(s) ds \end{aligned} \quad (3.9)$$

for  $0 < t, s \leq 1$ . The dual space of the space of continuous periodic functions is the space of periodic finite regular Borel measures equipped with the total variation norm and the duality formulation gives

$$\sup_{\substack{f: \|f\|_{\infty} \leq \frac{\sqrt{d}}{4} \\ \text{Lip}(f) \leq 1}} \min_{p \in \mathcal{P}^{n,d,\infty}} \|f - p\|_{\infty} = \sup_{\substack{f: \|f\|_{\infty} \leq \frac{\sqrt{d}}{4} \\ \text{Lip}(f) \leq 1}} \sup_{\substack{\hat{\mu}^{(k)=0, \|k\|_{\infty} \leq n \\ \|\mu\|_{\text{TV}} \leq 1}} \left| \int_{\mathbb{T}^d} f(x) d\mu(x) \right|.$$

Our main contribution to this result is the observation how to transfer the multivariate setting back to the univariate one. It is easy to verify that  $f(x) = \frac{1}{d^{3/2}} \sum_{\ell=1}^d f_0(x_{\ell})$  for a univariate Lipschitz function  $f_0$ ,  $\text{Lip}(f_0) \leq d$ ,  $\|f_0\|_{\infty} \leq \frac{d}{4}$  fulfils the conditions for the Lipschitz function  $f$ . Additionally,  $\mu^* = \frac{1}{d} \sum_{s=1}^d \mu_s$  with  $\mu_s = \left( \bigotimes_{\ell \neq s} \lambda(x_{\ell}) \right) \otimes \mu_0^*(x_s)$ ,

$$\mu_0^*(x_s) = \frac{1}{2(n+1)} \sum_{j=0}^{2n+1} (-1)^j \delta_{j/(2n+2)}(x_s)$$

and  $\lambda$  being the Lebesgue measure on  $\mathbb{T}$  is admissible<sup>57</sup>. Since this choice of  $\mu_s$  integrates  $\int g d\mu_s = 0$  if  $g$  is constant with respect to  $x_s$  (and the same holds for constant univariate

<sup>55</sup>One can easily see that the Fourier series of  $g(x) = \frac{1}{2} - x$ ,  $x \in [0, 1]$ , is given by the series in (3.8). By the Dirichlet-Jordan test, one directly obtains the convergence of the Fourier series towards  $g(x)$  at  $x \in (0, 1)$  and towards zero at the discontinuity point  $x = 0$ .

<sup>56</sup>This result is well-known as *Rademacher's theorem*, see for instance [141, Box 1.9].

<sup>57</sup>Note that  $\|\mu^*\|_{\text{TV}} \leq \frac{1}{d} \sum_{s=1}^d \|\lambda\|_{\text{TV}}^{d-1} \|\mu_0^*\|_{\text{TV}} = 1$  and

$$d \cdot \hat{\mu}^*(k) = \sum_{s=1}^d \hat{\mu}_0^*(k_s) \prod_{\ell \neq s} \delta_{k_{\ell}, 0} = \sum_{s=1}^d \sum_{j=0}^{2n+1} \frac{e^{-2\pi i j \frac{n+1+k_s}{2n+2}}}{2(n+1)} \prod_{\ell \neq s} \delta_{k_{\ell}, 0} = \sum_{s=1}^d \delta_{k_s, n+1+(2n+2)\mathbb{Z}} \prod_{\ell \neq s} \delta_{k_{\ell}, 0} = 0$$

for  $\|k\|_{\infty} \leq n$ . Within this calculation  $\delta_{i,j}$  for indices  $i, j \in \mathbb{Z}$  denotes the usual *Kronecker delta* being one if  $i = j$  and zero if  $i \neq j$ .

functions integrated against  $\mu_0^*$ ), we obtain with (3.9)

$$\begin{aligned}
 \sup_{\substack{f: \mathbb{T}^d \rightarrow \mathbb{R} \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}, \text{Lip}(f) \leq 1}} \min_{\substack{p \in \mathcal{P}^{n, d, \infty} \\ \hat{p}(0) = 1}} \|f - p\|_\infty &\geq \sup_{\substack{f_0: \mathbb{T} \rightarrow \mathbb{R} \\ \|f_0\|_\infty \leq \frac{d}{4}, \text{Lip}(f_0) \leq d}} \left| \frac{1}{d^{5/2}} \sum_{s, \ell=1}^d \int_{\mathbb{T}^d} f_0(x_\ell) d\mu_s(x) \right| \\
 &= \sup_{\substack{f_0: \mathbb{T} \rightarrow \mathbb{R} \\ \|f_0\|_\infty \leq \frac{d}{4}, \text{Lip}(f_0) \leq d}} \left| \frac{1}{d^{5/2}} \sum_{\ell=1}^d \int_{\mathbb{T}} f_0(x_\ell) d\mu_0^*(x_\ell) \right| \\
 &= \sup_{\substack{f_0: \mathbb{T} \rightarrow \mathbb{R} \\ \|f_0\|_\infty \leq \frac{d}{4}, \text{Lip}(f_0) \leq d}} \left| \int_{\mathbb{T}} \frac{f_0'(s)}{d^{3/2}} \left( \int_{\mathbb{T}} \mathcal{B}_1(t-s) d\mu_0^*(t) \right) ds \right|.
 \end{aligned}$$

We denote  $\mathcal{B}_{\mu^*}(s) = \int_{\mathbb{T}} \mathcal{B}_1(t-s) d\mu_0^*(t)$  and observe  $\int_{\mathbb{T}} \mathcal{B}_{\mu^*}(s) ds = 0$ . Moreover,  $\mu_0^*$  has moments  $\hat{\mu}_0^*(k) = 1$  for  $k \in (n+1)(2\mathbb{Z}+1)$  and  $\hat{\mu}_0^*(k) = 0$  otherwise. Together with the Fourier representation (3.8) of  $\mathcal{B}_1$  where one rewrites the sum over odd integers as the difference between the sum over all nonzero integers and the sum of all nonzero even integers, this gives

$$\begin{aligned}
 \mathcal{B}_{\mu^*}(s) &= \sum_{k \in \mathbb{Z} \setminus \{0\}} \frac{e^{2\pi i k(n+1)s}}{2\pi i k(n+1)} - \sum_{k \in \mathbb{Z} \setminus \{0\}} \frac{e^{2\pi i 2k(n+1)s}}{2\pi i 2k(n+1)} \\
 &= \frac{1}{n+1} \mathcal{B}_1((n+1)s) - \frac{1}{2n+2} \mathcal{B}_1((2n+2)s) \\
 &= \begin{cases} \frac{1}{4} \frac{1}{n+1}, & (n+1)s - \lfloor (n+1)s \rfloor \in (0, \frac{1}{2}), \\ -\frac{1}{4} \frac{1}{n+1}, & (n+1)s - \lfloor (n+1)s \rfloor \in (\frac{1}{2}, 1), \\ 0, & (2n+2)s \in \{0, \dots, 2n+1\}. \end{cases}
 \end{aligned}$$

Here, the last equality is a direct consequence of (3.8). Now, we choose  $f_0$  by taking  $f_0'(s) = d \cdot \text{sgn}(\mathcal{B}_{\mu^*}(s))$  and  $f_0(0) = 0$  which is possible as it yields

$$\|f_0\|_\infty = \frac{d}{2n+2} \leq \frac{d}{4} \quad \text{and} \quad \int_{\mathbb{T}} f_0'(s) ds = 0$$

for  $n \geq 1$ . Finally, we end up with

$$\begin{aligned}
 \sup_{\substack{f: \mathbb{T}^d \rightarrow \mathbb{R} \\ \|f\|_\infty \leq \frac{\sqrt{d}}{4}, \text{Lip}(f) \leq 1}} \min_{\substack{p \in \mathcal{P}^{n, d, \infty} \\ \hat{p}(0) = 1}} \|f - p\|_\infty &\geq d^{-1/2} \int_{\mathbb{T}} |\mathcal{B}_{\mu^*}(s)| ds \\
 &= \frac{d^{-1/2}}{n+1} \int_{\mathbb{T}} \left| \mathcal{B}_1((n+1)s) - \frac{1}{2} \mathcal{B}_1((2n+2)s) \right| ds \\
 &= \frac{d^{-1/2}}{n+1} \int_{\mathbb{T}} \left| \mathcal{B}_1(s) - \frac{1}{2} \mathcal{B}_1(2s) \right| ds = \frac{d^{-1/2}}{4(n+1)}. \quad \square
 \end{aligned}$$

We remark at this point that our lower bound on the best approximation differs from the one described in [26] by a factor  $d^{-1/2}$  due to the different definition of  $W_1$  which is based on  $|\cdot|_1$  as the underlying metric on  $\mathbb{T}^d$  in [26]. However, this does not change the statement too much as the lower bound is still sharp in the univariate case. The latter will be a consequence of the next subsection.

### 3 Trigonometric polynomials and rational functions

**Remark 3.1.7** (Information theoretic point of view). One should distinguish the above result on the best approximation by a polynomial with given degree  $n$  from the question of how well one can recover any measure given its low order trigonometric moments. While the polynomial approximation calculated in the framework of Theorem 3.1.6 is based on the knowledge of all moments, the latter information theoretic question would only consider the moments  $\hat{\mu}(k)$  for  $k \in \{-n, \dots, n\}^d$ . A lower bound can be reformulated as the largest difference

$$\sup \left\{ W_1(\mu, \nu) : \mu, \nu \in \mathcal{M}(\mathbb{T}^d), \hat{\nu}(k) = \hat{\mu}(k) \text{ for } k \in \{-n, \dots, n\}^d \right\} \quad (3.10)$$

between two measures, which have equal low order moments and cannot be distinguished by a recovery algorithm if no additional prior is known. If  $\hat{\mu}$  and  $\hat{\nu}$  are equal up to order  $n$ , then convolution with the Jackson kernel yields  $J_n * \mu = J_n * \nu$ , so that the triangle inequality for  $W_1$  and Remark 3.1.4 give

$$W_1(\mu, \nu) \leq W_1(\mu, J_n * \mu) + W_1(\nu, J_n * \nu) \leq \frac{3d}{2} \frac{\|\mu\|_{\text{TV}} + \|\nu\|_{\text{TV}}}{n+2},$$

and thus (3.10) is at most of order  $\mathcal{O}(n^{-1})$ . This order is also optimal which can be seen by choosing  $\mu$  as the Lebesgue measure  $\lambda$ ,  $\nu$  being absolutely continuous with  $d\nu(x_1, \dots, x_d) = [1 + \cos(2\pi(n+1)x_1)] d\lambda(x_1, \dots, x_d)$ , and  $f(x) = \cos(2\pi(n+1)x_1)/(2\pi(n+1))$  in

$$W_1(\mu, \nu) = \sup_{f: \text{Lip}(f) \leq 1} \int_{\mathbb{T}^d} f(x) \cos(2\pi(n+1)x_1) dx \geq \int_{\mathbb{T}^d} \frac{\cos^2(2\pi(n+1)x_1)}{2\pi(n+1)} dx = \frac{1}{4\pi(n+1)}.$$

This shows that the knowledge of the Fourier coefficients of a measure up to order  $n$  without any prior assumption on the measure only allows to approximate the measure with worst case error of order  $n^{-1}$ . This worst case error rate can be decreased if prior knowledge on the ground truth measure, e.g. sparsity as in Chapter 2, is assumed. There, we have even seen that  $\mu = \nu$  if two sufficiently nicely separated measures have equal moments  $\hat{\nu}(k) = \hat{\mu}(k)$  for  $k \in \{-n, \dots, n\}^d$ , see Theorem 2.2.8.

#### 3.1.4 Univariate situation and uniqueness of best approximation

On the univariate torus  $\mathbb{T}$ , the Wasserstein distance of two probability measures can be rewritten as a  $L^1$  distance of their cumulative density functions (CDF) shifted by some constant depending on the measures, see [21]. We extend this to real signed measures.

**Lemma 3.1.8** (Wasserstein via CDF). *For any univariate  $\mu, \nu \in \mathcal{M}_{\mathbb{R}}(\mathbb{T})$  with normalisation  $\mu(\mathbb{T}) = \nu(\mathbb{T}) = 1$ , we have*

$$W_1(\mu, \nu) = \int_0^1 |\mu([0, x]) - \nu([0, x]) - c^*(\mu, \nu)| dx,$$

and  $c^*(\mu, \nu) \in \mathbb{R}$  depends on  $\mu, \nu$ .

*Proof.* For  $\mu, \nu \in \mathcal{M}_{+,1}(\mathbb{T})$  this is [21, Thm. 3.7]. For signed  $\mu, \nu \in \mathcal{M}_{\mathbb{R}}(\mathbb{T})$ , we can use the Jordan decomposition of any signed measure as a difference of nonnegative measures. In other words, we write  $\mu = \mu_+ - \mu_-$ ,  $\nu = \nu_+ - \nu_-$  for  $\mu_+, \mu_-, \nu_+, \nu_- \in \mathcal{M}_+(\mathbb{T})$  and rewrite

$$W_1(\mu, \nu) = \sup_{f: \text{Lip}(f) \leq 1} \left| \int_{\mathbb{T}} f(x) [d\nu_+(x) + d\mu_-(x) - (d\nu_-(x) + d\mu_+(x))] \right|$$

$$= (\nu_+ + \mu_-)(\mathbb{T}) \cdot \sup_{f: \text{Lip}(f) \leq 1} \left| \int_{\mathbb{T}} f(x) \left[ \frac{d\nu_+(x) + d\mu_-(x)}{(\nu_+ + \mu_-)(\mathbb{T})} - \frac{d\nu_-(x) + d\mu_+(x)}{(\nu_+ + \mu_-)(\mathbb{T})} \right] \right| \quad (3.11)$$

and this allows to apply [21, Thm. 3.7].<sup>58</sup> This then gives

$$\begin{aligned} & \int_0^1 \left| \left( \frac{\nu_+ + \mu_- - (\nu_- + \mu_+)}{(\nu_+ + \mu_-)(\mathbb{T})} \right) ([0, x]) - c^* \left( \frac{\nu_+ + \mu_-}{(\nu_+ + \mu_-)(\mathbb{T})}, \frac{\nu_- + \mu_+}{(\nu_+ + \mu_-)(\mathbb{T})} \right) \right| dx \\ &= (\nu_+ + \mu_-)(\mathbb{T})^{-1} \int_0^1 |(\nu - \mu)([0, x]) - c^*(\nu, \mu)| dx \end{aligned}$$

for the Wasserstein distance of  $\mu$  and  $\nu$  where the constant  $c^*(\nu, \mu)$  depends again only on the two measures.  $\square$

As  $W_1$  is not strictly convex, the question of uniqueness of the best approximation with respect to this norm is not trivial. However, it can be equivalently characterised by the uniqueness of the best approximation in  $L^1(\mathbb{T})$  and thus allows for the following theorem.

**Theorem 3.1.9** (Best approximation in the univariate case). *If  $\mu, \nu \in \mathcal{M}_{\mathbb{R}}(\mathbb{T})$  with normalisation  $\mu(\mathbb{T}) = \nu(\mathbb{T}) = 1$  are absolutely continuous or give mass only to countably many atoms, we have*

$$W_1(\nu, \mu) = \inf_{c \in \mathbb{R}} \int_{\mathbb{T}} |(\mathcal{B}_1 * \nu)(t) - (\mathcal{B}_1 * \mu)(t) - c| dt. \quad (3.12)$$

*This allows to conclude that for any  $n \in \mathbb{N}$ , any real measure being normalised and absolutely continuous with respect to the Lebesgue measure admits a unique best approximation by a normalised polynomial of degree  $n \in \mathbb{N}$  with respect to the 1-Wasserstein distance.*

*Proof.* Let  $\mu, \nu \in \mathcal{M}_{\mathbb{R}}(\mathbb{T})$  and  $\mathcal{B}_1$  denote the Bernoulli spline of degree 1 from the proof of Theorem 3.1.6, then we have

$$\begin{aligned} W_1(\nu, \mu) &= \sup_{f: \text{Lip}(f) \leq 1} \left| \int_{\mathbb{T}} f(x) [d\nu(x) - d\mu(x)] \right| \\ &= \sup_{f: \text{Lip}(f) \leq 1} \left| \int_{\mathbb{T}} \int_{\mathbb{T}} f'(t) \mathcal{B}_1(x-t) [d\nu(x) - d\mu(x)] dt \right| \\ &= \sup_{f: \text{Lip}(f) \leq 1} \left| \int_{\mathbb{T}} f'(t) [(\mathcal{B}_1 * \nu)(t) - (\mathcal{B}_1 * \mu)(t)] dt \right|. \end{aligned}$$

Since the integral over  $f'$  is zero by the periodicity of  $f$ , any  $c \in \mathbb{R}$  yields

$$\begin{aligned} \left| \int_{\mathbb{T}} f'(t) [(\mathcal{B}_1 * \nu)(t) - (\mathcal{B}_1 * \mu)(t)] dt \right| &= \left| \int_{\mathbb{T}} f'(t) [(\mathcal{B}_1 * \nu)(t) - (\mathcal{B}_1 * \mu)(t) - c] dt \right| \\ &\leq \inf_{c \in \mathbb{R}} \int_{\mathbb{T}} |(\mathcal{B}_1 * \nu)(t) - (\mathcal{B}_1 * \mu)(t) - c| dt. \end{aligned}$$

We proceed by computing explicitly

$$(\mathcal{B}_1 * \mu)(t) = \int_{[0, t) \cup (t, 1)} \mathcal{B}_1(t-x) d\mu(x)$$

<sup>58</sup>Note that by  $0 = \nu(\mathbb{T}) - \mu(\mathbb{T}) = (\nu_+ + \mu_-)(\mathbb{T}) - (\mu_+ + \nu_-)(\mathbb{T})$  both measures in the integral in (3.11) are probability measures. Moreover, observe that  $(\nu_+ + \mu_-)(\mathbb{T}) \geq \nu(\mathbb{T}) = 1 > 0$  and hence it is possible to normalise as stated.

### 3 Trigonometric polynomials and rational functions

$$\begin{aligned}
&= \int_{[0,t)} \frac{1}{2} - (t-x) d\mu(x) + \int_{(t,1)} \frac{1}{2} - (t-x+1) d\mu(x) \\
&= \left(\frac{1}{2} - t\right) (\mu([0,1]) - \mu(\{t\})) + \int_{[0,1)} x d\mu(x) - t\mu(\{t\}) - \mu([0,1]) + \mu([0,t]) \\
&= \frac{\mu([0,t]) + \mu([0,t])}{2} - \mu([0,1]) \left(t + \frac{1}{2}\right) + \int_{[0,1)} x d\mu(x) \tag{3.13}
\end{aligned}$$

for  $t \in (0, 1)$  and

$$(\mathcal{B}_1 * \mu)(0) = \int_{[0,1)} x d\mu(x) - \frac{1}{2}\mu([0,1]) + \frac{1}{2}\mu(\{0\}).$$

On the other hand, Lemma 3.1.8 and (3.13) yield

$$\begin{aligned}
\int_0^1 |(\nu - \mu)([0, x]) - c^*(\nu, \mu)| dx &= W_1(\nu, \mu) \\
&\leq \inf_{c \in \mathbb{R}} \int_0^1 \left| (\nu - \mu)([0, x]) - \frac{(\nu - \mu)(\{x\})}{2} - c \right| dx
\end{aligned}$$

and thus equality (3.12) for absolutely continuous measures, or more generally for measures that give mass to at most countably many atoms, and such that the set of  $x$  where the integrands (in the equation above) from the upper and lower bounds disagree has Lebesgue measure zero.

With this knowledge, the question of approximation of  $\mu$  by  $p$  with degree  $n$  and  $\mathbf{p}_0 = 1$  can be rewritten as

$$\min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = \mu(0) = 1}} W_1(p, \mu) = \min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} \inf_{c \in \mathbb{R}} \int_{\mathbb{T}} |(\mathcal{B}_1 * p)(t) - (\mathcal{B}_1 * \mu)(t) - c| dt$$

if  $\mu$  does not give mass to single points. In this case, we have that  $\mathcal{B}_1 * \mu$  is continuous by (3.13), and hence there exists a unique best  $L^1$ -approximation  $\tilde{p} = \mathcal{B}_1 * p^* - c$  (see e.g. [36, Thm. 3.10.9]) which defines the unique best approximation  $p^*$  to  $\mu$  uniquely by  $\tilde{p} = \mathcal{B}_1 * p^* - c$  and the normalisation condition  $\mathbf{p}_0^* = 1$ .  $\square$

**Example 3.1.10.** Uniqueness and non-uniqueness of  $L^1$  approximation is discussed in some detail in [118, 41] and we note the following:

(i) For  $\mu = \frac{1}{2}\delta_0 - \frac{1}{2}\delta_{1/2} + \lambda \in \mathcal{M}_{\mathbb{R}}(\mathbb{T})$  where  $\lambda$  is again the Lebesgue measure, one finds

$$(\mathcal{B}_1 * \mu)(t) = \frac{1}{2} \left( \mathcal{B}_1(t) - \mathcal{B}_1\left(t - \frac{1}{2}\right) \right) = \begin{cases} 0, & t = 0, \\ \frac{1}{4}, & t \in (0, \frac{1}{2}), \\ 0, & t = \frac{1}{2}, \\ -\frac{1}{4}, & t \in (\frac{1}{2}, 1). \end{cases}$$

For any normalised polynomial  $p$ , we have that the difference  $\mathcal{B}_1 * p(t) - \mathcal{B}_1 * \mu(t)$  differs from  $\int_0^t p(x) dx - \mu([0, t])$  by a constant except at the discontinuity points  $t = 0, \frac{1}{2}$ . But as they have Lebesgue measure zero, we can derive from Theorem 3.1.9

$$\min_{\substack{p \in \mathcal{P}^{n,d,\infty} \\ \mathbf{p}_0 = 1}} W_1(p, \mu) = \inf_{c \in \mathbb{R}} \int_{\mathbb{T}} |(\mathcal{B}_1 * p)(t) - (\mathcal{B}_1 * \mu)(t) - c| dt. \tag{3.14}$$

As proven in [118, Thm. 5.1], the function  $\mathcal{B}_1 * \mu$  does not have a unique  $L^1$  approximation for even  $n$ . Thus,  $\mu$  does not admit a unique best approximation either.

- (ii) For  $\mu = \delta_0$  one has  $\mathcal{B}_1 * \mu = \mathcal{B}_1$  such that again (3.14) holds for this choice of  $\mu$ . According to [118, Lem. 2.2], this function  $\mathcal{B}_1$  with only one jump has a unique best  $L^1$ -approximation given by the interpolation polynomial

$$\tilde{p}(x) = \sum_{j=1}^n \frac{1}{2n+2} \cot\left(\frac{j\pi}{2n+2}\right) \sin(2\pi jx).$$

Deconvolving  $\tilde{p} = \mathcal{B}_1 * p^*$  gives

$$p^*(x) = 1 + \sum_{j=1}^n \frac{j\pi}{n+1} \cot\left(\frac{j\pi}{2n+2}\right) \cos(2\pi jx)$$

as the unique best approximation to  $\delta_0$ . Since the error of the best  $L^1$  approximation of  $\mathcal{B}_1$  is known from a theorem by Favard [51] (e.g. this is mentioned in [36, p. 213]), we can compute

$$\begin{aligned} W_1(\delta_0, p^*) &= \inf_{c \in \mathbb{R}} \int_{\mathbb{T}} |(\mathcal{B}_1 * p^*)(t) - (\mathcal{B}_1 * \delta_0)(t) - c| dt \\ &\leq \|\mathcal{B}_1 * p^* - \mathcal{B}_1\|_{L^1(\mathbb{T})} = \frac{1}{4(n+1)}. \end{aligned}$$

By comparison with Theorem 3.1.6, we notice that equality holds in this calculation and that the bound from Theorem 3.1.6 is sharp.

Figure 3.3 and Table 3.1 summarise our findings on the approximation of  $\delta_0$ . The best approximation  $p^*$  as well as the Dirichlet kernel  $D_n(x) = \sin((2n+1)\pi x)/\sin(\pi x)$  are signed with small full width at half maximum (FWHM) but positive and negative oscillations at the sides. The latter might be seen as an unwanted artifact in applications. The approximations given by the Fejér and the Jackson kernel are nonnegative. For

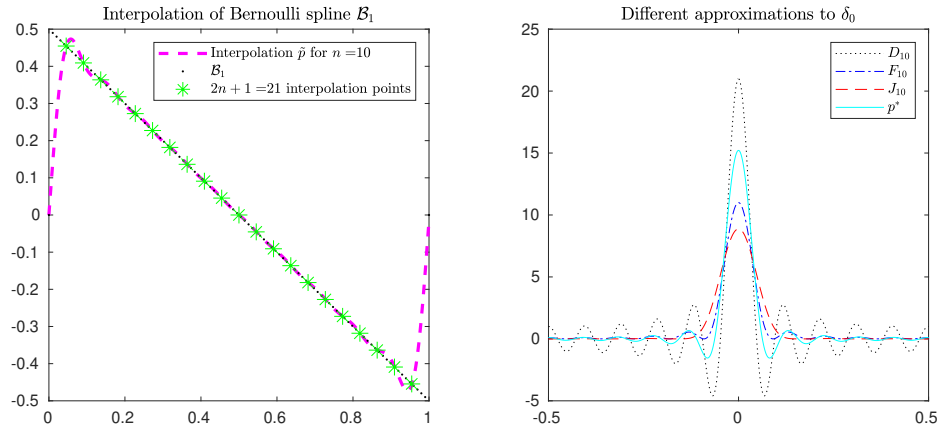


Figure 3.3: Interpolation of  $\mathcal{B}_1$  (left) and comparison of different polynomial approximations of degree  $n = 10$  to  $\delta_0$  (right).

completeness, we note that the Dirichlet kernel is the Fourier partial sum of  $\delta_0$  and allows for the estimate

$$W_1(\delta_0, D_n) \leq W_1(\delta_0, p^*) + W_1(p^*, D_n) \leq (1 + \|D_n\|_1) W_1(\delta_0, p^*) \leq \frac{\frac{4}{\pi^2} \log(n) + \mathcal{O}(1)}{4(n+1)} \quad (3.15)$$

### 3 Trigonometric polynomials and rational functions

which relies on  $W_1(p^*, D_n) = W_1(D_n * p^*, D_n * \delta_0) \leq \|D_n\|_1 W_1(\delta_0, p^*)$ , the well known bound on the Lebesgue constant [20, Prop. 1.2.3], and Example 3.1.10 (ii).<sup>59</sup>

Trig. polynomial	Sign of polynomial	$W_1(\delta_0, K_n)$
Dirichlet $D_n$	signed	$\leq \frac{\frac{4}{\pi^2} \log(n) + \mathcal{O}(1)}{4(n+1)}$ (Inequality (3.15))
Fejér $F_n$	nonnegative	$\leq \frac{1}{\pi^2} \frac{\log(n)+4}{n+1}$ (Theorem 3.1.3)
Jackson $J_n$ , $n$ even	nonnegative	$\leq \frac{3}{2} \frac{1}{n+2}$ (Remark 3.1.4)
Best approximation $p^*$	signed	$= \frac{1}{4(n+1)}$ (Example 3.1.10 (ii))

Table 3.1: Convergence rates of different trigonometric polynomials approximating the Dirac delta  $\delta_0$ .

**Remark 3.1.11.** We close by some remarks which are specific for the univariate setting:

- (i) We stress that Theorem 3.1.9 allows to compute the Wasserstein distance as an  $L^1$ -distance for real signed univariate measures. Similarly, this allows to compute the so-called star discrepancy  $\|\nu([0, \cdot])\|_\infty$  as suggested in [113, eq. (2.1) and (2.2)]. However note that (3.13) has some additional term such that  $\nu = \frac{1}{2}\delta_0 - \frac{1}{2}\delta_{1/2}$  with  $\nu(\mathbb{T}) = 0$  gives

$$\|\nu([0, \cdot])\|_\infty = \frac{1}{2} \neq \frac{1}{4} = \|\mathcal{B}_1 * \nu\|_\infty$$

and thus [113, eq. (2.1) and (2.2)] needs some adjustment. More precisely, it seems that in the publication [113] a factor  $\frac{1}{2}$  was lost since the  $k$ th Fourier coefficient of  $\nu([0, \cdot])$  is  $\frac{\hat{\nu}(k)}{ik}$  whereas  $\hat{\mathcal{B}}_1(k) \cdot \hat{\nu}(k) = \frac{1}{2} \frac{\hat{\nu}(k)}{ik}$ .

- (ii) In the univariate case, one can relate our work to a main result in [113]. As Theorem 3.1.9 reformulates the Wasserstein distance of two univariate measures in terms of the  $L^1$ -distance of their convolution with the Bernoulli spline, one can view this Bernoulli spline as a kernel of type  $\beta = 1$  following the notation of [113]. Thus, one can take  $p = 1, p' = \infty$  in [113, Thm. 4.1] yielding that the Wasserstein distance between a measure  $\mu$  and its trigonometric approximation is bounded from above by  $c/n$ . The latter agrees with our Remark 3.1.4 which additionally gives an explicit and small constant.
- (iii) The observation that the construction of  $p^*$  for  $\delta_0$  is possible via FFTs might lead to the idea to construct near-best approximations to any measure  $\mu$  by interpolating  $\mathcal{B}_1 * \mu$  by some  $\tilde{p}$  and to obtain the polynomial  $p$  of near best approximation which satisfies  $\tilde{p} = \mathcal{B}_1 * p$  by dividing with the Fourier coefficients of the Bernoulli spline  $\mathcal{B}_1$ . A first problem would be that the limited knowledge of moments only allows to interpolate the partial Fourier sum  $S_n(\mathcal{B}_1 * \mu)$ , which does not converge to  $\mathcal{B}_1 * \mu$  uniformly as  $n \rightarrow \infty$  for discrete  $\mu$ . Secondly, the near-best approximation  $p$  cannot be expected to be nonnegative for a nonnegative measure  $\mu$  which is another drawback compared to convolution with nonnegative kernels like the Fejér or Jackson kernel.

<sup>59</sup>Going through the lines of the proof of [20, Prop. 1.2.3], one could also find an explicit upper bound for the term hidden by  $\mathcal{O}(1)$  in (3.15).



## 3.2 Interpolation and approximation by the signal polynomial

While Section 3.1 focuses on weak approximations of a general measure  $\mu$ , in particular via convolution with smooth kernels, we consider in this section another type of polynomial estimator, denoted by  $p_{1,n}$  (3.19), which depends non-linearly on  $\mu$  and is able to identify at a finite degree its support, under some assumptions on the latter. More precisely, we restrict ourselves to discrete measures  $\mu \in \mathcal{M}(q)$  as in Chapter 2 and the main result of this section, stated in Theorem 3.2.6 below, is a quantitative rate for the pointwise convergence

$$p_{1,n}(x) \xrightarrow{n \rightarrow \infty} \mathbb{1}_{\text{supp } \mu}(x) = \begin{cases} 1, & x \in \text{supp } \mu, \\ 0, & \text{otherwise} \end{cases} \quad (3.16)$$

to the indicator function of the support. After discussing algebraic properties of this estimator (Subsection 3.2.1), we consider the case of discrete measures (Subsection 3.2.2) whereas measures with support on an algebraic variety were included in our publication [26]. For those, we could show a convergence similar to (3.16) but we do not describe this in detail at this point because this work focuses on sparse super resolution for discrete measures.

In contrast to [26] and Section 3.1, we come back to the radial setting where we take moments on a ball into account and which we studied already in Chapter 2 as well as in the introduction to this chapter. Therefore, let  $[n] := B_{n/2}(0) \cap \mathbb{Z}^d$  and  $N := |[n]|$ .<sup>60</sup> We use bold type to designate vectors (resp. matrices) of  $\mathbb{C}^N$  (resp.  $\mathbb{C}^{N \times N}$ ) only (vectors of  $\mathbb{T}^d$  or  $\mathbb{N}^d$  are left in normal type). We write

$$\mathbf{e}_x^{(n)} := (e^{-2\pi i k x})_{k \in [n]} \in \mathbb{C}^N$$

for the vector containing all  $d$ -variate trigonometric monomials up to Euclidean degree  $n/2$ . We often identify a polynomial  $p \in \langle e^{-2\pi i k \cdot}; k \in [n] \rangle$  with its vector of coefficients  $\mathbf{p} \in \mathbb{C}^N$ , i.e.

$$p(x) = \mathbf{e}_x^{(n)*} \mathbf{p} \quad \forall x \in \mathbb{T}^d.$$

Note that from Parseval's theorem,  $\|p\|_{L^2} = \|\mathbf{p}\|_2$ . Moreover, we highlight that  $|p|^2 \in \mathcal{P}^{n,d,2}$ . The key object of this section is the (truncated) moment matrix associated with the unknown measure  $\mu$ , defined as

$$\mathbf{T}_n := (\hat{\mu}(k - \ell))_{k, \ell \in [n]} \in \mathbb{C}^{N \times N}, \quad (3.17)$$

where  $\hat{\mu}(k)$  uses the trigonometric moments of  $\mu$  from (3.1). In this and the following section, we assume that these moments are exact whereas the last section of this chapter extends our analysis to noisy data.

### 3.2.1 Algebraic considerations

It is well-known that the range and kernel of the matrix (3.17) reveal some of the structure of the measure hidden behind the moments, and methods that aim at recovering  $\mu$  using purely algebraic manipulations on  $\mathbf{T}_n$  are often referred to as *subspace methods*, e.g. MUSIC

<sup>60</sup>Asymptotically, we thus have  $N \approx \text{Vol}_d(B_1(0))2^{-d}n^d$  where  $\text{Vol}_d(B_1(0))$  is the volume of the  $d$ -dimensional unit ball.

### 3 Trigonometric polynomials and rational functions

[144], ESPRIT [137] or matrix pencil [74]. The starting point for these methods is often the singular value decomposition of  $\mathbf{T}_n$ , see Section 1.1, which we denote by

$$\mathbf{T}_n = \mathbf{U}_n \mathbf{\Sigma}_n \mathbf{V}_n^* = \sum_{j=1}^N \sigma_j^{(n)} \mathbf{u}_j^{(n)} \mathbf{v}_j^{(n)*},$$

where all matrices are of size  $N \times N$ ,  $\mathbf{u}_j^{(n)}$  and  $\mathbf{v}_j^{(n)}$  are the  $j$ -th columns of  $\mathbf{U}_n$  and  $\mathbf{V}_n$  respectively (left and right singular vectors), and  $\sigma_1^{(n)} \geq \sigma_2^{(n)} \geq \dots \geq \sigma_n^{(n)}$  are the diagonal entries of the diagonal matrix  $\mathbf{\Sigma}_n$  (singular values). This decomposition is sometimes explicitly used to design estimators for the support of  $\mu$ , such as MUSICs frequency estimation function [144], or Christoffel polynomials [111]. In fact, it is interesting as a motivating remark to see that the construction of  $p_n = F_n * \mu$  from the previous section can also be expressed in terms of this singular value decomposition. As we changed our setting towards the radial data, one has to highlight that the classical Fejér kernel  $F_n(x) = \frac{|d_n(x)|^2}{(n+1)^d}$  is replaced by its radial analogue  $F_{\text{rad},n}(x) := \frac{|D_{\text{rad},n/2}(x)|^2}{N}$  compared to [26, Lem. 4.1]. Here,  $D_{\text{rad},n/2}$  is the radial Dirichlet kernel defined in Section 1.3.

**Lemma 3.2.1.** *The moment matrix  $\mathbf{T}_n$  fulfils*

$$(F_{\text{rad},n} * \mu)(x) = \frac{1}{N} \mathbf{e}_x^{(n)*} \mathbf{T}_n \mathbf{e}_x^{(n)} = \frac{1}{N} \sum_{j=1}^N \sigma_j^{(n)} u_j^{(n)}(x) \overline{v_j^{(n)}(x)}, \quad (3.18)$$

where, as explained above,  $u_j^{(n)}(x) = \mathbf{e}_x^{(n)*} \mathbf{u}_j^{(n)}$  and  $v_j^{(n)}(x) = \mathbf{e}_x^{(n)*} \mathbf{v}_j^{(n)}$ .

*Proof.* We have for any  $x \in \mathbb{T}^d$

$$\begin{aligned} \frac{1}{N} \mathbf{e}_x^{(n)*} \mathbf{T}_n \mathbf{e}_x^{(n)} &= \frac{1}{N} \sum_{k \in [n]} \sum_{l \in [n]} \hat{\mu}(k-l) e^{2\pi i k x} e^{-2\pi i l x} \\ &= \frac{1}{N} \int_{\mathbb{T}^d} \sum_{k \in [n]} \sum_{l \in [n]} e^{-2\pi i (k-l)y} e^{2\pi i k x} e^{-2\pi i l x} d\mu(y) \\ &= \int_{\mathbb{T}^d} \frac{1}{N} \left| \sum_{k \in [n]} e^{-2\pi i k (y-x)} \right|^2 d\mu(y) = \int_{\mathbb{T}^d} F_{\text{rad},n}(x-y) d\mu(y), \end{aligned}$$

where the last equality is by definition of  $F_{\text{rad},n}$  as a generalisation of (3.3). Plugging in the singular value decomposition of  $\mathbf{T}_n$  yields the second equality of the statement.  $\square$

Note that if  $\mu \in \mathcal{M}_{\mathbb{R}}(\mathbb{T}^d)$  (the set of real-valued measures), then the moment matrix  $\mathbf{T}_n$  is Hermitian. If  $\mu \in \mathcal{M}_+(\mathbb{T}^d)$  (the set of nonnegative measures), then  $\mathbf{T}_n$  is positive semi-definite, and we have in particular the sum of squares representation

$$(F_{\text{rad},n} * \mu)(x) = \frac{1}{N} \sum_{j=1}^N \sigma_j^{(n)} \left| v_j^{(n)}(x) \right|^2.$$

We now introduce polynomial estimators for the measure, which can be understood as the *unweighted* counterparts of  $p_n$ . Let  $r_n := \text{rank } \mathbf{T}_n$  and define *signal- and noise-polynomials*  $p_{1,n}, p_{0,n} : \mathbb{T}^d \rightarrow [0, 1]$  respectively, by

$$p_{1,n}(x) = \frac{1}{N} \sum_{j=1}^{r_n} \left| v_j^{(n)}(x) \right|^2 \quad \text{and} \quad p_{0,n}(x) = \frac{1}{N} \sum_{j=r_n+1}^N \left| v_j^{(n)}(x) \right|^2. \quad (3.19)$$

### 3.2 Interpolation and approximation by the signal polynomial

This signal/noise terminology comes from the notions of signal and noise subspaces, which were initially introduced in [144] and are at the core of the aforementioned subspace methods in signal processing (we refer the interested reader to [110, Section 9.6] for an overview). Schematically speaking, they correspond to the spaces spanned by the vectors  $(\mathbf{v}_1^{(n)}, \dots, \mathbf{v}_{r_n}^{(n)})$  (the *signal* space) and  $(\mathbf{v}_{r_n+1}^{(n)}, \dots, \mathbf{v}_N^{(n)})$  (the *noise* space) respectively. They are actually independent of the singular value decomposition itself, which ensures in particular that  $p_{1,n}$  and  $p_{0,n}$  are indeed well-defined.

The key idea of subspace methods, relating these spaces to the underlying measure  $\mu$ , is that, given a polynomial  $p \in \langle e^{-2\pi i k x}; k \in [n] \rangle$  that vanishes on  $\text{supp } \mu$ , one obtains using (3.17) that the  $k$ -th entry ( $k \in [n]$ ) of  $\mathbf{T}_n \mathbf{p}$  is given by

$$\sum_{l \in [n]} \mathbf{p}_l \cdot \int_{\mathbb{T}^d} e^{-2\pi i (k-l)x} d\mu(x) = \int_{\mathbb{T}^d} e^{-2\pi i k x} p(x) d\mu(x) = 0, \quad (3.20)$$

and thus  $\mathbf{p} \in \ker \mathbf{T}_n$ . Hence, finding the common roots of all polynomials contained in the kernel of the matrix  $\mathbf{T}_n$  may allow to identify the support of  $\mu$ . In what follows, we denote by  $V(\ker \mathbf{T}_n)$  the variety consisting of the common roots of all the polynomials in  $\ker \mathbf{T}_n$ , i.e.

$$V(\ker \mathbf{T}_n) := \{x \in \mathbb{T}^d : p(x) = \mathbf{e}_x^{(n)*} \mathbf{p} = 0 \text{ for all } \mathbf{p} \in \ker \mathbf{T}_n\}.$$

We begin in this section with qualitative, purely algebraic considerations about the polynomials (3.19). The next theorem shows that, under the condition that  $\text{supp } \mu = V(\ker \mathbf{T}_n)$ ,  $p_{0,n}$  and  $p_{1,n}$  actually identify the support of  $\mu \in \mathcal{M}(q)$  for *finite*  $n$ .

**Theorem 3.2.2.** *Let  $d, n \in \mathbb{N}, q > 0$ ,  $\mu \in \mathcal{M}(q)$ , and suppose  $V(\ker \mathbf{T}_n) = \text{supp } \mu \subseteq \mathbb{T}^d$ . Then  $p_{0,n}(x) + p_{1,n}(x) = 1$  for all  $x \in \mathbb{T}^d$ . In particular, we have*

$$p_{1,n}(x) \begin{cases} = 1, & \text{if } x \in \text{supp } \mu, \\ < 1, & \text{otherwise.} \end{cases} \quad (3.21)$$

*Proof.* We have

$$p_{1,n}(x) + p_{0,n}(x) = \frac{1}{N} \sum_{j=1}^N \left| v_j^{(n)}(x) \right|^2 = \frac{1}{N} \mathbf{e}_x^{(n)*} \mathbf{V}_n \mathbf{V}_n^* \mathbf{e}_x^{(n)} = \frac{1}{N} \mathbf{e}_x^{(n)*} \mathbf{e}_x^{(n)} = 1, \quad (3.22)$$

so in particular  $p_{1,n}(x) \in [0, 1]$ . Since  $V(\ker \mathbf{T}_n) = \text{supp } \mu$  and  $\ker \mathbf{T}_n = \langle \mathbf{v}_{r_n+1}^{(n)}, \dots, \mathbf{v}_N^{(n)} \rangle$ , it follows that the polynomials  $v_{r_n+1}^{(n)}, \dots, v_N^{(n)}$  vanish on  $\text{supp } \mu$ , so  $p_{1,n}(x) = 1$  for all  $x \in \text{supp } \mu$ . Conversely, if  $x \in \mathbb{T}^d$  such that  $p_{1,n}(x) = 1$ , we claim that  $x \in \text{supp } \mu$ . Indeed, we have  $1 - p_{1,n}(x) = \sum_{j=r_n+1}^N \left| v_j^{(n)}(x) \right|^2 = 0$ , so it follows that  $x$  lies in the vanishing set of  $v_{r_n+1}^{(n)}, \dots, v_N^{(n)}$ , so  $x \in V(\ker \mathbf{T}_n) = \text{supp } \mu$ .  $\square$

**Remark 3.2.3.** For more general measures  $\mu$  with support on an algebraic variety we have extended Theorem 3.2.2 in [26]. Variants of this result are known for discrete measures and for measures with support on curves or surfaces, e.g. in [86] and [124, Propositions 5.2, 5.3], respectively. The hypothesis  $V(\ker \mathbf{T}_n) = \text{supp } \mu$  in Theorem 3.2.2 is well-known in the theory of super resolution [89, 142] or polynomial system solving [94], and is hard to check in practice. For discrete measures, this sufficient condition can be guaranteed if the geometry of the support behaves sufficiently nicely, i.e. if the support points are well

### 3 Trigonometric polynomials and rational functions

separated from each other. Interestingly, we derive in Theorem 3.2.6 that a separation of the nodes slightly larger than the Rayleigh condition which we introduced in Chapter 2 as a necessary condition for a small condition number of the problem is also sufficient in order to guarantee  $V(\ker \mathbf{T}_n) = \text{supp } \mu$ .

**Example 3.2.4.** For  $\mu = \delta_0$ , we have  $p_{1,n}(x) = N^{-1} \cdot F_{\text{rad},n}(x)$  by Lemma 3.2.1 and the proofs of the Theorems 3.2.6 and 3.2.9 will also show that  $p_{1,n}$  is generally close to a sum of normalised Fejér kernels for well-separated discrete measures.

We conclude this subsection by stating a variational characterisation of  $p_{0,n}$ , which will be a useful tool in proofs within the next subsections.

**Lemma 3.2.5.** *If  $\ker \mathbf{T}_n \neq \{\mathbf{0}\}$ , we have that*

$$p_{0,n}(x) = \max \left\{ \frac{1}{N} \frac{|p(x)|^2}{\|\mathbf{p}\|_2^2} : \mathbf{p} \in \ker \mathbf{T}_n \setminus \{\mathbf{0}\} \right\}. \quad (3.23)$$

*Proof.* As we assume  $\ker \mathbf{T}_n \neq \{\mathbf{0}\}$ , we have  $r_n := \text{rank } \mathbf{T}_n < N$  and find a rectangular matrix  $\mathbf{V}_0 = (\mathbf{v}_{r_n+1}^{(n)}, \dots, \mathbf{v}_N^{(n)}) \in \mathbb{C}^{N \times (N-r_n)}$  whose columns form an orthonormal basis of  $\ker \mathbf{T}_n$ . For fixed  $x \in \mathbb{T}^d$ , let  $\mathbf{q}_x := \mathbf{V}_0 \mathbf{V}_0^* \mathbf{e}_x^{(n)} \in \ker \mathbf{T}_n$  such that we identify this vector of coefficients with the polynomial satisfying

$$\mathbf{q}_x(x) = \mathbf{e}_x^{(n)*} \mathbf{q}_x = \sum_{j=r_n+1}^N \left| v_j^{(n)}(x) \right|^2 = N p_{0,n}(x).$$

For all  $\mathbf{p} \in \ker \mathbf{T}_n$ , we have

$$\mathbf{q}_x^* \mathbf{p} = \mathbf{e}_x^{(n)*} \mathbf{V}_0 \mathbf{V}_0^* \mathbf{p} = \mathbf{e}_x^{(n)*} \mathbf{p} = p(x).$$

In particular, note that

$$\|\mathbf{q}_x\|_2^2 = \mathbf{q}_x^* \mathbf{q}_x = \mathbf{e}_x^{(n)*} \mathbf{V}_0 \mathbf{V}_0^* \mathbf{e}_x^{(n)} = N p_{0,n}(x). \quad (3.24)$$

Therefore, by the Cauchy–Schwarz inequality, it follows that

$$|p(x)|^2 = |\mathbf{q}_x^* \mathbf{p}|^2 \leq \|\mathbf{q}_x\|_2^2 \cdot \|\mathbf{p}\|_2^2 = N p_{0,n}(x) \cdot \|\mathbf{p}\|_2^2.$$

Hence, we have

$$p_{0,n}(x) \geq \max_{\mathbf{p} \in \ker \mathbf{T}_n \setminus \{\mathbf{0}\}} \frac{|p(x)|^2}{N \|\mathbf{p}\|_2^2} \geq \frac{|q_x(x)|^2}{N \|\mathbf{q}_x\|_2^2} = p_{0,n}(x),$$

if  $\mathbf{q}_x \neq \mathbf{0}$ . The first inequality also holds when  $\mathbf{q}_x = \mathbf{0}$ , in which case the result follows due to (3.24).  $\square$

### 3.2.2 The signal polynomial for discrete measures

We now come to the first main result of this section, stated in Theorem 3.2.6 below, which gives quantitative rates for the pointwise convergence (3.16) in the case where  $\mu$  is a discrete measure. If the measure is given by

$$\mu = \sum_{t \in Y} \alpha_t \delta_t$$

with support  $\text{supp } \mu = Y \subset \mathbb{T}^d$  and complex, nonzero weights  $\alpha_t$  forming the matrix  $W = \text{diag}((\alpha_t)_{t \in Y})$ , then the moment matrix allows for the *Vandermonde factorisation*

$$\mathbf{T}_n = \mathbf{A}_n^* W \mathbf{A}_n, \quad \mathbf{A}_n = (e^{2\pi i k t})_{\substack{t \in Y \\ k \in [n]}} \in \mathbb{C}^{|Y| \times N}.$$

### 3.2 Interpolation and approximation by the signal polynomial

**Theorem 3.2.6** (Pointwise convergence). *Let  $\mu = \sum_{t \in Y} \alpha_t \delta_t$ ,  $\alpha_t \in \mathbb{C}$ ,  $Y \subset \mathbb{T}^d$ , and let  $x \in \mathbb{T}^d$  such that  $x \neq t$  for all  $t \in Y$ . Let  $\alpha_{\min}$  and  $\alpha_{\max}$  be the minimal and maximal weights in absolute value. If  $(n-1) \cdot \text{sep } Y = \frac{2\sqrt{1+\tau}j_{d/2-1,1}}{\pi}$  for some  $\tau > 0$ , there is  $c_{d,\tau}^{(8)} > 0$  such that*

$$p_{1,n}(x) \leq \frac{\alpha_{\max}}{\alpha_{\min}} \cdot \frac{c_{d,\tau}^{(8)}}{n^2} \sum_{t \in Y} \frac{1}{\|x - t\|_{\mathbb{T}^d}^2}$$

for  $x \notin Y$ . In particular, this implies the pointwise convergence (3.16). Under the stronger assumption  $(n-1) \cdot \text{sep } Y = \frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi}$  for some  $\tau > 0$ , one additionally finds a constant  $c_{d,\tau}^{(9)} > 0$  such that for  $x \in \mathbb{T}^d$  with  $\min_{t \in Y} \|x - t\|_{\mathbb{T}^d} \leq \frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi n}$  we have

$$p_{1,n}(x) \leq 1 - c_{d,\tau}^{(9)} n^2 \min_{t \in Y} \|x - t\|_{\mathbb{T}^d}^2.$$

*Proof.* The condition  $(n-1) \cdot \text{sep } Y = \frac{2\sqrt{1+\tau}j_{d/2-1,1}}{\pi}$  implies  $\text{rank } \mathbf{A}_{n-1} = |Y| < N$  due to work on minorant functions in [59].<sup>61</sup> Since the weights are nonzero and  $\text{rank } \mathbf{A}_{n-1} \leq \text{rank } \mathbf{A}_n \leq |Y|$ , this yields  $\text{rank } \mathbf{T}_n = \text{rank } \mathbf{A}_n = |Y|$ . Based on this, [86, Thm. 2.8] gives  $\text{supp } \mu = V(\ker \mathbf{A}_n) = V(\ker \mathbf{T}_n)$  such that this separation condition allows to apply Theorem 3.2.2.<sup>62</sup> We have  $p_{1,n} = 1 - p_{n,0}$  by Theorem 3.2.2 and because of  $\langle \mathbf{v}_{|Y|+1}, \dots, \mathbf{v}_N \rangle = \ker \mathbf{T}_n = \ker \mathbf{A}_n$ , it follows that the space  $\langle \mathbf{v}_1, \dots, \mathbf{v}_{|Y|} \rangle$  does not depend on the weights  $\alpha_t$ . Thus, we can consider the absolute value measure  $|\mu| = \sum_{t \in Y} |\alpha_t| \delta_t$  and know that the range of the moment matrix  $\widetilde{\mathbf{T}}_n = \widetilde{\mathbf{V}}_n \widetilde{\Sigma}_n \widetilde{\mathbf{V}}_n^*$  corresponding to  $|\mu|$  agrees with the range of  $\mathbf{T}_n$ . Therefore, one has for any  $x \in \mathbb{T}^d$

$$p_{1,n}(x) = \frac{1}{N} \sum_{j=1}^{|Y|} |v_j^{(n)}(x)|^2 = \frac{1}{N} \sum_{j=1}^{|Y|} |\tilde{v}_j^{(n)}(x)|^2 \leq \frac{1}{N} \sum_{j=1}^{|Y|} \frac{\tilde{\sigma}_j^{(n)}}{\sigma_{\min}(\widetilde{\mathbf{T}}_n)} |\tilde{v}_j^{(n)}(x)|^2$$

and the application of (3.18) together with Lemma 1.3.7 gives

$$p_{1,n}(x) \leq \frac{(F_{\text{rad},n} * |\mu|)(x)}{\sigma_{\min}(\widetilde{\mathbf{T}}_n)} = \frac{\sum_{t \in Y} |\alpha_t| F_{\text{rad},n}(x-t)}{\sigma_{\min}(\widetilde{\mathbf{T}}_n)} \leq \frac{\sqrt{d} c_d^2 \alpha_{\max} (n/2)^{2d-2}}{N \sigma_{\min}(\widetilde{\mathbf{T}}_n)} \sum_{t \in Y} \frac{4}{\|x - t\|_{\mathbb{T}^d}^2}.$$

Then, the first statement follows by observing that  $N \in \mathcal{O}(n^d)$  and<sup>63</sup>

$$\sigma_{\min}(\widetilde{\mathbf{T}}_n) = \min_{\mathbf{u} \neq 0} \frac{\mathbf{u}^* \mathbf{A}_n^* |W| \mathbf{A}_n \mathbf{u}}{\|\mathbf{u}\|_2^2} \geq \alpha_{\min} \sigma_{\min}(\mathbf{A}_n)^2 \geq \alpha_{\min} c_{d,\tau}^{(7)} n^d$$

for some  $c_{d,\tau}^{(7)} > 0$  by Proposition 2.3.2.

For the second part, we denote the  $(|Y|+1)$ -th standard basis vector by  $e_{|Y|+1} = (0, \dots, 0, 1)^\top \in \mathbb{C}^{|Y|+1}$  and consider the Vandermonde matrix

$$\tilde{\mathbf{A}}_{n,x} = \begin{bmatrix} \mathbf{A}_n^* & \mathbf{e}_x^{(n)} \end{bmatrix} \in \mathbb{C}^{N \times (|Y|+1)}.$$

<sup>61</sup>The work by Goncalves [59] shows that there exists a minorant function  $\hat{\psi}$  for  $\mathbb{1}_{B_{(n-1)/2}(0)}$  whose inverse Fourier transform  $\psi$  is supported in  $B_{\frac{2\sqrt{1+\tau}j_{d/2-1,1}}{\pi(n-1)}}(0)$  and admits  $\psi(0) > 0$ . This implies  $\text{rank } \mathbf{A}_{n-1} =$

$\min(|Y|, N) = |Y|$ , see [86, Thm. 2.4 and Cor. 2.5].

<sup>62</sup>Even though [86, Thm. 2.8] is formulated for the case of max-degree frequencies  $k \in \{0, \dots, n\}^d$  instead of the radial frequency sampling considered in this section, one can directly translate their proof for  $\text{supp } \mu = V(\ker \mathbf{A}_n)$  if  $\text{rank } \mathbf{A}_{n-1} = |Y|$  to the radial setting.

<sup>63</sup>We write  $|W|$  if we apply the absolute value entrywise to  $W$ .

### 3 Trigonometric polynomials and rational functions

We want to obtain a bound on  $p_{0,n}$  by plugging an admissible test polynomial into (3.23). The pseudo-inverse of  $\tilde{\mathbf{A}}_{n,x}$  gives rise to a candidate given by the Lagrange polynomial

$$\ell_{|Y|+1}(y) = e_{|Y|+1}^* \tilde{\mathbf{A}}_{n,x}^\dagger \mathbf{e}_y^{(n)},$$

satisfying  $\ell_{|Y|+1}(t) = 0$  for  $t \in Y$  and  $\ell_{|Y|+1}(x) = 1$ .<sup>64</sup> We compute

$$\begin{aligned} \|\ell_{|Y|+1}\|_{L^2}^2 &= \int_{\mathbb{T}^d} |e_{|Y|+1}^* \tilde{\mathbf{A}}_{n,x}^\dagger \mathbf{e}_y^{(n)}|^2 dy \\ &= \int_{\mathbb{T}^d} |\langle \tilde{\mathbf{A}}_{n,x}^{\dagger*} e_{|Y|+1}, \mathbf{e}_y^{(n)} \rangle|^2 dy = \|\tilde{\mathbf{A}}_{n,x}^{\dagger*} e_{|Y|+1}\|_2^2 \leq \sigma_{\min}(\tilde{\mathbf{A}}_{n,x})^{-2} \end{aligned}$$

and use Lemma 3.2.5 to bound

$$1 - p_{1,n}(x) = p_{0,n}(x) = \max_p \frac{|p(x)|^2}{N \|p\|_{L^2}^2} \geq \frac{|\ell_{|Y|+1}(x)|^2}{N \|\ell_{|Y|+1}\|_{L^2}^2} \geq \frac{\sigma_{\min}(\tilde{\mathbf{A}}_{n,x})^2}{N}.$$

Finally, this allows to conclude the assertion by using Theorem 2.2.8.<sup>65</sup>  $\square$

The above theorem shows that  $p_{1,n}$  converges pointwisely to the indicator function of the support of the measure which we want to recover if the support points are separated by more than  $\frac{2j_{d/2-1,1}}{\pi(n-1)}$ . On the contrary, we have seen in Chapter 2 that a separation of at least  $\frac{j_{d/2,1}}{\pi n}$  is sufficient for the possible existence of a stable algorithm. For  $d = 1$ , the sufficient condition for the algorithm and the necessary condition for well-conditionedness are asymptotically equal to  $n^{-1}$  as  $n \rightarrow \infty$  whereas in the bivariate case the sufficient condition  $\frac{2j_{0,1}}{\pi(n-1)} \approx 1.53(n-1)^{-1}$  is only slightly larger than the Rayleigh length being approximately  $1.22n^{-1}$  and the necessary condition  $\sqrt{\frac{4}{3}}n^{-1}$ .

**Remark 3.2.7.** Actually, Theorem 3.2.6 shows the correct orders in  $n$  and  $\min_{t \in Y} |x - t|_\infty^2$  in the upper bound of  $p_{1,n}(x)$ . First note that  $1 - p_{1,n}$  and all its partial derivatives of order 1 vanish on  $Y$ . For fixed  $x \in \mathbb{T}^d$ , and  $t' = \operatorname{argmin}_{t \in Y} |x - t|_\infty$ , the Taylor expansion at  $t'$  thus gives  $\xi \in \mathbb{T}^d$  such that

$$1 - p_{1,n}(x) = \frac{1}{2} (x - t')^\top H_x(\xi) (x - t'),$$

where  $H_x(\xi) := (-\partial_s \partial_t p_{1,n}(\xi))_{1 \leq s, t \leq d}$  is the Hessian of  $1 - p_{1,n}$  at  $\xi$ . Thus,

$$1 - p_{1,n}(x) \leq \frac{1}{2} \|H_x(\xi)\|_F \cdot |x - t'|_2^2 \leq \frac{d}{2} \max_{r,s} \|\partial_r \partial_s p_{1,n}\|_{L^\infty} \cdot d |x - t'|_\infty^2.$$

<sup>64</sup>Recall that  $\tilde{\mathbf{A}}_{n,x}$  has full rank due to the separation of the nodes, see Theorem 2.2.8 or Footnote

65. Hence, we have  $\ell_{|Y|+1}(t) = e_{|Y|+1}^* \tilde{\mathbf{A}}_{n,x}^\dagger \mathbf{e}_t^{(n)} = e_{|Y|+1}^* \tilde{\mathbf{A}}_{n,x}^\dagger \tilde{\mathbf{A}}_{n,x} e_t = e_{|Y|+1}^* e_t = 0$  for  $t \in Y$  and analogously  $\ell_{|Y|+1}(x) = e_{|Y|+1}^* e_{|Y|+1} = 1$ .

<sup>65</sup>The variational formulation of the smallest singular value gives

$$\sigma_{\min}(\tilde{\mathbf{A}}_{n,x})^2 = \min_{\substack{(v,v')^\top \in \mathbb{C}^{|Y|+1} \\ \|(v,v')^\top\|_2=1}} \|\mathbf{A}_n^* v + \mathbf{e}_x^{(n)} v'\|_2^2$$

and the latter can be analysed easily by the lower bound on the distance of two moments vectors corresponding to well-separated measures which we presented in Theorem 2.2.8.

### 3.2 Interpolation and approximation by the signal polynomial

One may apply Bernstein's inequality (see e.g. [36, Chapter 4]) to  $y_s \mapsto p_{1,n}(y_1, \dots, y_d)$  and  $y_r \mapsto \partial_s p_{1,n}(y_1, \dots, y_d)$  successively (both univariate trigonometric polynomials of degree  $n$ ), and obtain

$$1 - p_{1,n}(x) \leq 2\pi^2 d^2 n^2 \cdot \min_{t \in Y} |x - t|_\infty^2$$

since  $\|p_{1,n}\|_{L^\infty} = 1$ . A bivariate visualisation of the bounds on  $p_{1,n}$  is shown in Figure 3.4.

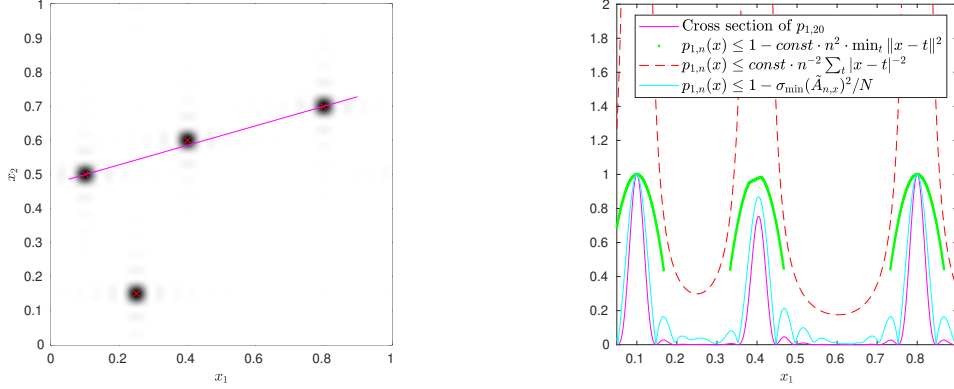


Figure 3.4: Summary of the bounds on  $p_{1,n}$  from Theorem 3.2.6 and Remark 3.2.7 for  $d = 2$ ,  $n = 20$ , and a discrete measure  $\mu$  supported on four points. The polynomial  $p_{1,20}$  was evaluated on a grid in  $\mathbb{T}^2$  and interpolated on the magenta cross section (left), while the bounds on  $p_{1,20}$  on this cross section are displayed (right). We see that specifically the bound  $1 - \sigma_{\min}(\tilde{\mathbf{A}}_{n,x})^2/N$  from the proof of Theorem 3.2.6 reproduces the behaviour of  $p_{1,n}$ .

In fact, normalising  $p_{1,n}$  differently even leads to a weak convergence result towards the empirical measure associated with the support points. This result, stated in Theorem 3.2.9 below, uses the following technical lemma.

**Lemma 3.2.8** (Convergence of singular values). *For  $Y = \{t_j : j = 1, \dots, |Y|\} \subset \mathbb{T}^d$  let  $\mu = \sum_{j=1}^{|Y|} \alpha_{t_j} \delta_{t_j}$  be a discrete complex measure whose weights are ordered non-increasingly with respect to their absolute value. Assume that  $n \cdot \text{sep } Y > \frac{2\sqrt{1+\tau}j_{d/2-1,1}}{\pi}$ , then there is a constant  $c_{d,\tau}^{(10)} > 0$  such that the singular values  $\sigma_j^{(n)}$ ,  $j = 1, \dots, N$ , of the moment matrix  $\mathbf{T}_n$  fulfil*

$$\left| |\alpha_{t_j}| - \frac{\sigma_j^{(n)}}{N} \right| \leq \frac{c_{d,\tau}^{(10)}}{n} \cdot \frac{\alpha_{\max}|Y|}{\text{sep } Y}, \quad j = 1, \dots, |Y|.$$

*Proof.* With the polar decomposition  $\frac{1}{\sqrt{N}} \mathbf{A}_n^* = \mathbf{P}\mathbf{H}$ , where  $\mathbf{P} \in \mathbb{C}^{N \times |Y|}$  is unitary and  $\mathbf{H} \in \mathbb{C}^{|Y| \times |Y|}$  is positive-definite, we have that  $|\alpha_{t_1}| \geq \dots \geq |\alpha_{t_{|Y|}}|$  are the singular values of the matrix  $\mathbf{P}\mathbf{W}\mathbf{P}^*$ . Therefore, we obtain for the singular values of  $\mathbf{T}_n = \mathbf{A}_n^* \mathbf{W} \mathbf{A}_n$

$$\begin{aligned} \max_{1 \leq j \leq |Y|} \left| \frac{\sigma_j^{(n)}}{N} - |\alpha_{t_j}| \right| &\leq \left\| \frac{1}{N} \mathbf{T}_n - \mathbf{P}\mathbf{W}\mathbf{P}^* \right\|_2 = \|\mathbf{H}\mathbf{W}\mathbf{H}^* - \mathbf{W}\|_2 \\ &\leq \|\mathbf{H}\mathbf{W}(\mathbf{H} - \mathbf{I}_{|Y|})\|_2 + \|(\mathbf{H} - \mathbf{I}_{|Y|})\mathbf{W}\|_2 \end{aligned}$$

### 3 Trigonometric polynomials and rational functions

$$\begin{aligned}
&\leq |\alpha_{\max}| (\|\mathbf{H}\|_2 + 1) \|\mathbf{H} - \mathbf{I}_{|Y|}\|_2 \\
&\leq |\alpha_{\max}| (\|\mathbf{H}\|_2 + 1) \|(\mathbf{H} + \mathbf{I}_{|Y|})^{-1}\|_2 \|\mathbf{H}^2 - \mathbf{I}_{|Y|}\|_2 \\
&\leq |\alpha_{\max}| \frac{\frac{1}{\sqrt{N}} \sigma_{\max}(\mathbf{A}_n) + 1}{\frac{1}{\sqrt{N}} \sigma_{\min}(\mathbf{A}_n) + 1} \left\| \frac{1}{N} \mathbf{A}_n \mathbf{A}_n^* - \mathbf{I}_{|Y|} \right\|_2,
\end{aligned}$$

where the first inequality is due to [15, Theorem 2.2.8] and the last inequality is a consequence of  $\mathbf{H} = \mathbf{P}^* \mathbf{P} \mathbf{H} = \frac{1}{\sqrt{N}} \mathbf{P}^* \mathbf{A}_n^*$  and

$$\mathbf{H}^2 = \mathbf{H}^* \mathbf{H} = \frac{1}{N} \mathbf{A}_n \mathbf{P} \mathbf{P}^* \mathbf{A}_n^* = \frac{1}{\sqrt{N}} \mathbf{A}_n \mathbf{P} \mathbf{P}^* \mathbf{P} \mathbf{H} = \mathbf{A}_n \frac{1}{\sqrt{N}} \mathbf{P} \mathbf{H} = \frac{1}{N} \mathbf{A}_n \mathbf{A}_n^*.$$

Each entry of the matrix  $\frac{1}{N} \mathbf{A}_n \mathbf{A}_n^* - \mathbf{I}_{|Y|}$  except for the main diagonal is a radial Dirichlet kernel such that Lemma 1.3.7 together with Gerschgorin's theorem give the uniform bound

$$\left\| \frac{1}{N} \mathbf{A}_n \mathbf{A}_n^* - \mathbf{I}_{|Y|} \right\|_2 = \frac{1}{N} \max_{j=1, \dots, |Y|} \sum_{l \neq j} \left| \sum_{k \in [n]} e^{2\pi i k(t_l - t_j)} \right| \leq \frac{|Y| - 1}{N} \cdot \frac{c_d \sqrt{d} (n/2)^{d-1}}{\text{sep } Y}.$$

Since we can note that  $\sigma_{\min}(\mathbf{A}_n) > 0$  and

$$\frac{1}{\sqrt{N}} \sigma_{\max}(\mathbf{A}_n) = \sqrt{\left\| \frac{1}{N} \mathbf{A}_n \mathbf{A}_n^* \right\|_2} \leq 1 + \sqrt{\left\| \frac{1}{N} \mathbf{A}_n \mathbf{A}_n^* - \mathbf{I}_{|Y|} \right\|_2},$$

we obtain the proposed result for some constant  $c_{d,\tau}^{(10)}$ .  $\square$

**Theorem 3.2.9.** *We have*

$$\frac{p_{1,n}}{\|p_{1,n}\|_{L^1}} \rightarrow \tilde{\mu} = \frac{1}{|Y|} \sum_{t \in Y} \delta_t$$

as  $n \rightarrow \infty$ .

*Proof.* First note that  $\|p_{1,n}\|_{L^1} = \frac{|Y|}{N}$  if  $n$  is sufficiently large such that  $\mathbf{A}_n$  has full rank  $|Y|$ . Our idea is to estimate the 1-Wasserstein distance between  $\frac{p_{1,n}}{\|p_{1,n}\|_{L^1}}$  and  $\tilde{\mu}$ . For this, we define  $\tilde{p}_n = F_{\text{rad},n} * \tilde{\mu}$  and observe that for any Lipschitz continuous function  $f$  on  $\mathbb{T}^d$  with  $\text{Lip}(f) \leq 1$ ,  $\|f\|_\infty \leq \frac{\sqrt{d}}{4}$ , we have

$$\begin{aligned}
&\left| \int_{\mathbb{T}^d} \frac{p_{1,n}(x)}{\|p_{1,n}\|_{L^1}} f(x) dx - \frac{1}{|Y|} \sum_{t \in Y} f(t) \right| \\
&\leq \left| \int_{\mathbb{T}^d} \left( \frac{p_{1,n}(x)}{\|p_{1,n}\|_{L^1}} - \tilde{p}_n(x) \right) f(x) dx \right| + \left| \int_{\mathbb{T}^d} \tilde{p}_n(x) f(x) dx - \frac{1}{|Y|} \sum_{t \in Y} f(t) \right| \\
&\leq \left\| \frac{N}{r} p_{1,n} - \tilde{p}_n \right\|_{L^1} \|f\|_{L^\infty} + \left| \int_{\mathbb{T}^d} f d(F_{\text{rad},n} * \tilde{\mu}) - \int_{\mathbb{T}^d} f d\tilde{\mu} \right| \\
&\leq \frac{\sqrt{d}}{4} \left\| \frac{N}{|Y|} p_{1,n} - \tilde{p}_n \right\|_{L^1} + \int_{\mathbb{T}^d} F_{\text{rad},n}(x) \|x\|_{\mathbb{T}^d} dx.
\end{aligned} \tag{3.25}$$

Hence, it is enough to show that  $\left\| \frac{N}{|Y|} p_{1,n} - \tilde{p}_n \right\|_{L^1}$  and the integral  $\int_{\mathbb{T}^d} F_{\text{rad},n}(x) \|x\|_{\mathbb{T}^d} dx$  converge to zero for  $n \rightarrow \infty$ .



### 3.2 Interpolation and approximation by the signal polynomial

If  $\bar{n}$  is sufficiently large, then by Lemma 3.2.1 we can write  $\tilde{p}_n(x) = \frac{1}{N} e_x^{(n)*} \tilde{U} \tilde{\Sigma} \tilde{U}^* e_x^{(n)}$ , where  $\tilde{\Sigma} \in \mathbb{C}^{|Y| \times |Y|}$  denotes the diagonal matrix consisting of non-zero singular values, and  $\tilde{U} \in \mathbb{C}^{N \times |Y|}$  denotes the corresponding singular vector matrix of the moment matrix of  $\tilde{\mu}$ . As  $p_{1,n}$  only depends on the signal space of the moment matrix  $\mathbf{T}_n$  of  $\mu$ , which agrees with the signal space of the moment matrix of  $\tilde{\mu}$ , it follows by (3.19) that  $p_{1,n}(x) = \frac{1}{N} e_x^{(n)*} \tilde{U} \tilde{U}^* e_x^{(n)}$  and thus

$$\left| \frac{N}{|Y|} p_{1,n}(x) - \tilde{p}_n(x) \right| = \left| e_x^{(n)*} \tilde{U} \left( \frac{\mathbf{I}_{|Y|}}{|Y|} - \frac{\tilde{\Sigma}}{N} \right) \tilde{U}^* e_x^{(n)} \right| \leq \left\| e_x^{(n)*} \tilde{U} \right\|_2^2 \left\| \frac{1}{|Y|} \mathbf{I}_{|Y|} - \frac{1}{N} \tilde{\Sigma} \right\|_2.$$

Because  $\int_{\mathbb{T}^d} \left\| e_x^{(n)*} \tilde{U} \right\|_2^2 dx = N \|p_{1,n}\|_{L^1} = |Y|$  is constant, the convergence

$$\left\| \frac{N}{|Y|} p_{1,n} - \tilde{p}_n \right\|_{L^1} \rightarrow 0$$

follows from Lemma 3.2.8. For the second term in (3.25), we denote the coefficients of  $F_{\text{rad},n}(x)$  by  $\mathbf{c}_k$  and because of  $F_{\text{rad},n}(x) = \frac{1}{N} |D_{\text{rad},n/2}(x)|^2$  they can be computed via the convolution

$$\mathbf{c}_{k \cdot e_1} = \frac{1}{N} \sum_{\substack{j \in \mathbb{Z}^d \\ \|j\|_2 \leq \frac{n}{2}, \|j - k \cdot e_1\|_2 \leq \frac{n}{2}}} 1 = 1 - \frac{1}{N} \sum_{\substack{j \in \mathbb{Z}^d \\ \|j\|_2 \leq \frac{n}{2}, \|j - k \cdot e_1\|_2 > \frac{n}{2}}} 1 \geq 1 - \frac{1}{N} \sum_{\substack{j \in \mathbb{Z}^d \\ \frac{n}{2} < \|j - k \cdot e_1\|_2 \leq \frac{n}{2} + k}} 1$$

for  $k \in \{1, 2, \dots, n\}$ .<sup>66</sup> Analogously to the proof of Lemma 3.1.2, we can then bound

$$\begin{aligned} \int_{\mathbb{T}^d} F_{\text{rad},n}(x) \|x\|_{\mathbb{T}^d} dx &\leq \sum_{s=1}^d \int_{\mathbb{T}^d} \sum_{\substack{k \in \mathbb{Z}^d \\ \|k\|_2 \leq n}} \mathbf{c}_k e^{2\pi i k x} |x_s|_{\mathbb{T}} dx \\ &= d \int_{\mathbb{T}} \sum_{k=-n}^n \mathbf{c}_{k \cdot e_1} e^{2\pi i k x} |x|_{\mathbb{T}} dx \\ &= d \left( \frac{1}{4} + \sum_{k=1}^n \mathbf{c}_{k \cdot e_1} \int_0^{\frac{1}{2}} 4 \cos(kx) x dx \right) \\ &= d \left( \sum_{\ell=\lfloor \frac{n-1}{2} \rfloor + 1}^{\infty} \frac{1}{\pi^2 (2\ell + 1)^2} + \sum_{\ell=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{2(1 - \mathbf{c}_{(2\ell+1)e_1})}{\pi^2 (2\ell + 1)^2} \right). \end{aligned}$$

We know already from the proof of Lemma 3.1.2 that the first sum admits a rate of  $\mathcal{O}(n^{-1})$  as  $n \rightarrow \infty$ . The second sum can be controlled if we use Taylor's theorem in the last step of the estimation

$$1 - \mathbf{c}_{k e_1} \leq \frac{1}{N} \sum_{\substack{j \in \mathbb{Z}^d \\ \frac{n}{2} < \|j - k \cdot e_1\|_2 \leq \frac{n}{2} + k}} 1 \leq \frac{2^d c_d}{n^d} \int_{\frac{n}{2}}^{\frac{n}{2} + k} r^{d-1} dr = \frac{c_d}{d} \left( \left(1 + \frac{2k}{n}\right)^d - 1 \right) \leq \frac{2c_d 2^d k}{dn}$$

<sup>66</sup>By definition of  $N$ , we have  $\mathbf{c}_0 = 1$ . Moreover, the symmetry property  $\mathbf{c}_{-k \cdot e_1} = \mathbf{c}_{k \cdot e_1} = \mathbf{c}_{k \cdot e_s}$  holds for all  $k \in \mathbb{Z}$  and all  $s \in \{1, 2, \dots, d\}$ .

### 3 Trigonometric polynomials and rational functions

for some constant  $c_d > 0$  in order to obtain

$$\sum_{\ell=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{2(1 - c_{(2\ell+1)e_1})}{\pi^2(2\ell+1)^2} \leq \mathcal{O} \left( n^{-1} \sum_{\ell=0}^{\lfloor \frac{n-1}{2} \rfloor} \frac{1}{2\ell+1} \right) = \mathcal{O} \left( \frac{\log(n)}{n} \right)$$

and this completes the proof.  $\square$

#### 3.2.3 Numerical examples for approximation and interpolation

We complete this section by an illustration of the asymptotic behaviour of  $p_n$  and  $p_{1,n}$  for discrete measures with respect to the 1-Wasserstein distance. We compute the distance using a semidiscrete optimal transport algorithm described below. The code to reproduce the figures is available at <https://github.com/Paulcat/Measure-trigo-approximations> and was implemented by Paul Catala [26]. The only difference is that we use the Euclidean metric on the torus in contrast to the 1-norm considered in [26] such that this results in a difference in the computation of the Wasserstein distance which is based on these norms on  $\mathbb{T}^2$ . The updated code can be found in the GitHub repository containing the relevant code of this work, see <https://github.com/MHockmann/Dissertation.git>.

Among the measures studied in the numerical examples of [26], we only present a bivariate, discrete measure  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  supported on 15 points with (nonnegative) random amplitudes because the focus of this work is not on measures with support on curves. The polynomials  $p_n$ ,  $J_n * \mu$ , and  $p_{1,n}$  can be evaluated efficiently via the fast Fourier transform over a regular grid in  $\mathbb{T}^2$ . For the polynomial  $p_{1,n}$ , the singular value decomposition of the moment matrix  $\mathbf{T}_n$  can be computed at reduced cost by exploiting that  $\mathbf{T}_n$  has Toeplitz structure and resorting only to matrix-vector multiplications which can be computed by means of the FFT. This improvement will play a more important role when we deal with more realistic data sets in Section 4.1.

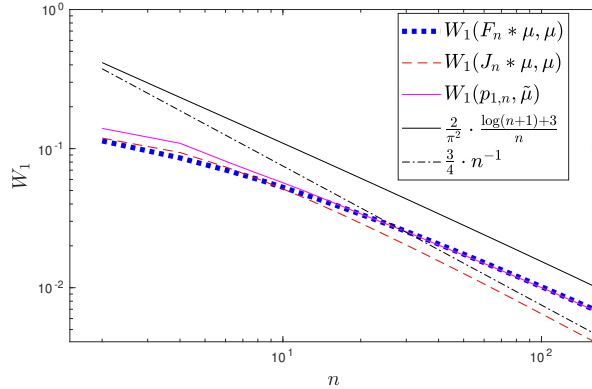


Figure 3.5: Asymptotics of  $p_n$  and  $p_{1,n}$ . For  $p_{1,n}$ , the distance is computed with respect to the unweighted measure  $\tilde{\mu} = \frac{1}{|Y|} \sum_{t \in Y} \delta_t$  where  $Y$  is the support of  $\mu$ .

The *semidiscrete optimal transport* between a measure with density  $p$  and the discrete measure  $\mu$  may be computed by solving the finite-dimensional optimisation problem

$$\max_{w \in \mathbb{R}_+^{|Y|}} f(w), \quad f(w) = \sum_{t \in Y} \alpha_t w_t + \sum_{t \in Y} \int_{\Omega_t(w)} (\|t - y\|_2 - \alpha_t) p(y) dy, \quad (3.26)$$

where the *Laguerre cells* associated to the weight vector  $w$  are given by

$$\Omega_t(w) = \left\{ y \in \mathbb{T}^d : \|t - y\|_2 - w_t \leq \|t' - y\|_2 - w_{t'} \text{ for all } t' \in Y \right\},$$

see e.g. [129]. In our implementation, the density measure (and the Laguerre cells) are computed over a  $641 \times 641$  grid. The maximisation (3.26) is performed by a BFGS algorithm using the Matlab implementation [143] and the iteration is stopped when the change of the value of the objective goes below  $10^{-9}$ , or when the norm  $\|\nabla f\|_\infty$  goes below  $10^{-5}$ . Note that this last condition has a geometrical interpretation since the  $j$ -th component of  $\nabla f$  corresponds to the difference between the measure of the Laguerre cell  $\Omega_{t_j}(w)$  and the amplitude  $\alpha_{t_j}$ . We set the limit number of iterations to 100.

Our numerical results depicted in Figure 3.5 show that the distance  $W_1(p_n, \mu)$  decreases at a rate close to the worst-case bound derived in Theorem 3.1.3. However, the multiplicative constant  $\frac{d}{\pi^2}$  might not be optimal as already suggested by the lower bound in Theorem 3.1.3. This is also the case for  $W_1(p_{1,n}, \tilde{\mu})$ , which is coherent with the bound given in the proof of Theorem 3.2.9. Additionally, we can observe the faster rate for convolution with the Jackson kernel.

### 3.3 Rational Christoffel functions

The polynomials  $u_j^{(n)}(x), v_j^{(n)}(x)$ ,  $j = 1, \dots, N$ , arising from the singular value decomposition of the moment matrix  $\mathbf{T}_n$  satisfy

$$\begin{aligned} \langle u_j^{(n)}, u_{j'}^{(n)} \rangle &= \langle v_j^{(n)}, v_{j'}^{(n)} \rangle = \int_{\mathbb{T}^d} \overline{v_j^{(n)}(x)} v_{j'}^{(n)}(x) dx = \langle \mathbf{v}_j^{(n)}, \mathbf{v}_{j'}^{(n)} \rangle = \delta_{j,j'} \text{ and} \\ \langle u_j^{(n)}, v_{j'}^{(n)} \rangle_\mu &:= \int_{\mathbb{T}^d} \overline{u_j^{(n)}(x)} v_{j'}^{(n)}(x) d\mu(x) = \mathbf{u}_j^{(n)*} \mathbf{T}_n \mathbf{v}_{j'}^{(n)} = \sigma_j^{(n)} \delta_{j,j'} \end{aligned} \quad (3.27)$$

such that they might be seen as *orthogonal polynomials*. The moment matrix  $\mathbf{T}_n$  is Hermitian and positive semidefinite if the underlying measure is nonnegative such that its SVD coincides with the eigendecomposition and  $\mathbf{u}_j^{(n)} = \mathbf{v}_j^{(n)}$  for all  $j = 1, \dots, \text{rank } \mathbf{T}_n$  can be assumed in this case. The interpolant  $p_{1,n}$  is then the sum of the squares of the first  $r_n = \text{rank } \mathbf{T}_n$  polynomials  $u_j^{(n)}$  being orthogonal to each other with respect to the Lebesgue measure and to the underlying measure  $\mu$ . However, these polynomials  $u_j^{(n)}$  differ slightly from the usual definition of orthogonal polynomials. First, note that in contrast to most of the theory of orthogonal polynomials, we consider the case of a measure that is discrete while orthogonal polynomials are usually considered with respect to an absolutely continuous measure. Secondly, the degree of the polynomial is independent of their index  $j = 1, \dots, N$ , as opposed to usual orthogonal polynomials having increasing degree. Finally, one should note that any sequence of polynomials orthogonal with respect to a measure supported on finitely many points is finite. Nevertheless, the connection to orthogonal polynomials directly motivates to study the *Christoffel function* known from the field of orthogonal polynomials, see [123] and the references therein for a summary or the more recent works for instance in [93, 128, 111, 85, 84]. In our setting from the previous section, where we construct the moment matrix  $\mathbf{T}_n$  of a measure  $\mu$  by  $\mathbf{T}_n := (\hat{\mu}(k - \ell))_{k, \ell \in [n]} \in \mathbb{C}^{N \times N}$  with  $[n] := B_{n/2}(0) \cap \mathbb{Z}^d$ , the Christoffel function can be defined as follows.

**Definition 3.3.1.** (Christoffel function, e.g. [123, 93, 128]) Let  $\mu \in \mathcal{M}_+(\mathbb{T}^d)$  and assume that its moments with respect to the vector of monomials  $\mathbf{e}_x^{(n)}$  are given. If the moment

### 3 Trigonometric polynomials and rational functions

matrix  $\mathbf{T}_n$  has full rank, the *Christoffel-Darboux polynomial* is given by  $\mathbf{e}_x^{(n)*} \mathbf{T}_n^{-1} \mathbf{e}_x^{(n)}$  and the reciprocal rational function

$$q_\mu : \mathbb{T}^d \rightarrow \mathbb{R}, \quad x \mapsto q_\mu(x) := \frac{1}{\mathbf{e}_x^{(n)*} \mathbf{T}_n^{-1} \mathbf{e}_x^{(n)}}$$

is called *Christoffel function*.

**Remark 3.3.2.** (i) The assumption on the rank of  $\mathbf{T}_n$  being equivalent to  $\ker \mathbf{T}_n = 0$  is fulfilled if  $\mu$  is a measure  $\mu \in \mathcal{M}_+(\mathbb{T}^d)$  whose support has a nonempty interior, cf. [111]. This is because the eigenvalues of  $\mathbf{T}_n$  are equal to the singular values for nonnegative measures and these satisfy (3.27) leading to

$$\sigma_j^{(n)} = \int_{\mathbb{T}^d} |u_j^{(n)}(x)|^2 d\mu(x) > 0$$

in this case such that  $\mathbf{T}_n$  is positive definite and in particular invertible. On the other hand, discrete measures typically lead to  $\text{rank } \mathbf{T}_n < N$  if  $n$  is large enough as we have seen in Section 3.2. Hence, this definition of the Christoffel function is not directly applicable to discrete measures.

(ii) The Christoffel-Darboux polynomial can be seen as an evaluation of the *Christoffel-Darboux kernel*  $k_{\mu,n}(x, y) := \mathbf{e}_x^{(n)*} \mathbf{T}_n^{-1} \mathbf{e}_y^{(n)}$ ,  $x, y \in \mathbb{T}^d$  for  $x = y$ . The kernel admits the *reproducing kernel property*

$$\int_{\mathbb{T}^d} k_{\mu,n}(x, y) p(y) d\mu(y) = \mathbf{e}_x^{(n)*} \mathbf{T}_n^{-1} \left( \int_{\mathbb{T}^d} \mathbf{e}_y^{(n)} \mathbf{e}_y^{(n)*} d\mu(y) \right) \mathbf{p} = \mathbf{e}_x^{(n)*} \mathbf{p} = p(x)$$

for any polynomial  $p \in \mathcal{P}^{n,d,2}$ , cf. [93, 128].

(iii) Using the SVD of the moment matrix, we have

$$q_\mu(x) = \frac{1}{\sum_{j=1}^N |u_j^{(n)}(x)|^2 / \sigma_j^{(n)}}.$$

If  $\text{rank } \mathbf{T}_n = r_n < N$ , it is natural, see [128, 111], to extend this by setting

$$q_\mu(x) = \begin{cases} \frac{1}{\sum_{j=1}^{r_n} |u_j^{(n)}(x)|^2 / \sigma_j^{(n)}} = \frac{1}{\mathbf{e}_x^{(n)*} \mathbf{T}_n^\dagger \mathbf{e}_x^{(n)}}, & \text{if } \text{proj}_{\ker \mathbf{T}_n}(\mathbf{e}_x^{(n)}) = 0, \\ 0, & \text{else.} \end{cases}$$

(iv) If we are given exact moments of a discrete measure  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  with  $\alpha_t > 0$  and  $\text{sep } Y$  large enough such that the corresponding Vandermonde matrix  $\mathbf{A}_n$  satisfies  $|Y| = \text{rank } \mathbf{A}_n = \text{rank } \mathbf{T}_n$  and  $V(\ker \mathbf{T}_n) = Y$ , we can conclude that  $\text{proj}_{\ker \mathbf{T}_n}(\mathbf{e}_x^{(n)}) = \text{proj}_{\mathcal{R}(\mathbf{A}_n)^\perp}(\mathbf{e}_x^{(n)})$  vanishes if and only if  $x \in Y$ . Therefore, we have under these assumptions by Lemma 1.1.8

$$q_\mu(x) = \begin{cases} \frac{1}{\mathbf{e}_x^{(n)*} \mathbf{A}_n^\dagger W^{-1} \mathbf{A}_n \mathbf{e}_x^{(n)}} = \alpha_t, & \text{if } x = t \in Y, \\ 0, & \text{else.} \end{cases}$$

A generalisation of this result can be found in [128, p. 249]. This type of Christoffel function gives an exact representation of the measure in theory but this is not useful in applications due to the sensitivity of the unregularised  $q_\mu$  to noise in the moments, e.g. cf. [111].

A well-known property of the Christoffel function  $q_\mu$  is that it can be linked to the underlying measure  $\mu$  in the following variational way.

**Lemma 3.3.3.** (Variational characterisation of  $q_\mu$ , e.g. [128, Lem. 1]) For  $\mu \in \mathcal{M}_+(\mathbb{T}^d)$ , we have

$$q_\mu(x) = \min \left\{ \int_{\mathbb{T}^d} |p(y)|^2 d\mu(y) : p \in \mathcal{P}^{n/2, d, 2}, p(x) = 1 \right\}$$

for any  $x \in \mathbb{T}^d$ . The minimiser  $\tilde{p}$  for the right hand side is  $\tilde{p}(y) = \frac{k_{\mu, n}(x, y)}{k_{\mu, n}(x, x)} \in \mathcal{P}^{n/2, d, 2}$ .

For measures with a density with respect to the Lebesgue measure, there exist results on the pointwise or uniform convergence of the Christoffel function to the density function, see [93, 128]. In contrast to this, we are interested in the recovery of sparse, discrete measures where we have argued in Remark 3.3.2 that regularisation would be needed in the presence of noise. Different regularisation schemes based on the singular values are proposed in [111] and among those we choose to regularise the Christoffel function for discrete measures with the following “spectral cut-off” approach. This comes with the benefit that we do not have to choose an additional parameter for an approximation of the rank of the moment matrix.

**Definition 3.3.4.** (Regularised Christoffel function, cf. [111]) Let  $\sigma_{|Y|}^{(n)} > 0$ . For any  $\mu \in \mathcal{M}(q) \cap \mathcal{M}_+(\mathbb{T}^d)$  and  $\varepsilon \in (0, \sigma_{|Y|}^{(n)})$ , we define for any  $x \in \mathbb{T}^d$

$$q_{\varepsilon, n}(x) := \left( \sum_{j=1}^N g_\varepsilon(\sigma_j^{(n)}) |u_j^{(n)}(x)|^2 \right)^{-1} = \frac{1}{\mathbf{e}_x^{(n)*} \mathbf{T}_n^\dagger \mathbf{e}_x^{(n)} + \frac{N}{\varepsilon} (1 - p_{1, n}(x))}$$

as the *regularised Christoffel function* where the spectral cut-off function is

$$g_\varepsilon : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_+, \quad \sigma_j^{(n)} \mapsto \begin{cases} \frac{1}{\sigma_j^{(n)}} & , \sigma_j^{(n)} > \varepsilon, \\ \frac{1}{\varepsilon} & , \sigma_j^{(n)} \leq \varepsilon. \end{cases}$$

In contrast to a Tikhonov regularisation approach using the positive definite moment matrix  $\mathbf{T}_n + \varepsilon \mathbf{I}_N$  corresponding to the measure  $\mu + \varepsilon \lambda$  where  $\lambda$  is the Lebesgue measure on  $\mathbb{T}^d$ , the cut-off regularised Christoffel function does not inherit a simple variational formulation as in Lemma 3.3.3. But we have in analogy to Lemma 3.3.3 that

$$q_{\varepsilon, n}(x) = \min \left\{ \mathbf{p}^* (\mathbf{T}_n + \varepsilon \mathbf{U}_{0, n} \mathbf{U}_{0, n}^*) \mathbf{p} : p \in \mathcal{P}^{n/2, d, 2}, p(x) = 1 \right\}, \quad (3.28)$$

where  $\mathbf{U}_{0, n} \in \mathbb{C}^{N \times (N - |Y|)}$  contains in its columns the singular vectors of  $\mathbf{T}_n$  which span the kernel of  $\mathbf{T}_n$ . Typically, upper bounds on Christoffel functions are obtained by plugging test polynomials into (3.28), see for instance [128, 93, 111]. In particular, the *needle polynomial* established in [85] is often used as a feasible polynomial  $p \in \mathcal{P}^{n/2, d, 2}$  with  $p(x) = 1$  and exponential decay away from  $x$ .

Due to the lack of such a simple variational formulation, we will use a different approach and derive bounds on  $q_{\varepsilon, n}$  by taking the SVD of  $\mathbf{T}_n$  into account because the previous section already gives an intuition how the singular values and the singular functions behave if  $\mu$  is a discrete measure. Exemplarily, the Christoffel function for  $\mu \in \mathcal{M}(q) \cap \mathcal{M}(\mathbb{T})$  is displayed in Figure 3.6 where for comparison a scaled version of the convolution with

### 3 Trigonometric polynomials and rational functions

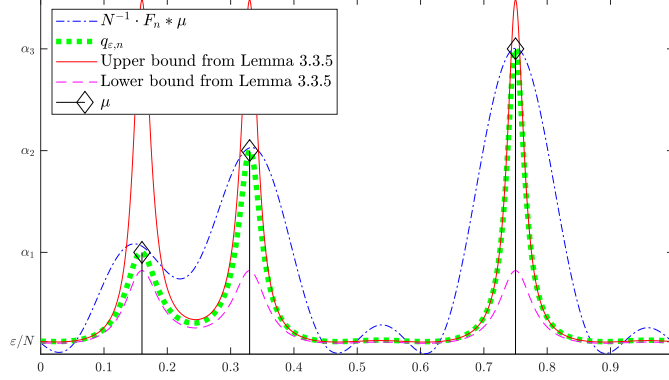


Figure 3.6: Christoffel function for a univariate, discrete measure  $\mu$  with three nodes and  $n = 6$ . In contrast to the convolution with the Fejér kernel as discussed in Section 3.1 (blue dash-dotted), the Christoffel function  $q_{\varepsilon,n}$  with  $\varepsilon = N^{-1}$  (green dotted), its upper bound derived in Lemma 3.3.5 (red solid) and the lower bound from Lemma 3.3.5 (magenta dashed) are more concentrated around the support of the ground truth measure.

$F_n$  as studied in Section 3.1 is included as well. A first observation is the interpolation property  $q_{\varepsilon,n}(t) = \alpha_t$  for all  $t \in Y$  which holds for any allowed positive value of  $\varepsilon$  if  $n$  is large enough such that  $\text{rank } \mathbf{A}_n = |Y|$ . Secondly, it follows from the definition that

$$q_{\varepsilon,n}(x) = \frac{1}{\frac{N}{\varepsilon} - \sum_{j=1}^{|Y|} \left( \frac{1}{\varepsilon} - \frac{1}{\sigma_j^{(n)}} \right) |u_j^{(n)}(x)|^2} \geq \frac{\varepsilon}{N}.$$

Thirdly, one can notice the following bounds in the plot of the Christoffel function.<sup>67</sup>

**Lemma 3.3.5** (Simple bounds on  $q_{\varepsilon,n}$ ). *Let  $\mu$  be a probability measure  $\mu \in \mathcal{M}(q) \cap \mathcal{M}_+(\mathbb{T}^d)$ ,  $\sigma_{|Y|}^{(n)} > 0$  and  $0 < \varepsilon < \sigma_{|Y|}^{(n)}$ . Then, we have for every  $x \in \mathbb{T}^d$*

$$\frac{\varepsilon}{N(1 - p_{1,n}(x)) + \frac{N\varepsilon}{\sigma_{|Y|}^{(n)}}} \leq q_{\varepsilon,n}(x) = \frac{\varepsilon}{N(1 - p_{1,n}(x)) + \varepsilon \mathbf{e}_x^{(n)*} \mathbf{T}_n^\dagger \mathbf{e}_x^{(n)}} \leq \frac{\varepsilon}{N(1 - p_{1,n}(x)) + \frac{N\varepsilon}{\sigma_1^{(n)}}}.$$

*Proof.* The upper bound is a direct consequence of  $\|\mathbf{e}_x^{(n)}\|_2^2 = N$  and the fact that the smallest eigenvalue value of  $\mathbf{T}_n^\dagger$  is  $1/\sigma_1^{(n)}$ . The lower bound follows analogously.  $\square$

This lemma yields that  $\lim_{\varepsilon \rightarrow 0} q_{\varepsilon,n}(x) = 0$  if  $p_{1,n}(x) < 1$  and by (3.21) we can conclude that this is the case exactly if  $x \notin \text{supp } \mu$  and  $n$  is sufficiently large in order to guarantee  $V(\ker \mathbf{T}_n) = \text{supp } \mu$ . Therefore, we already observe that  $q_{\varepsilon,n}$  has the ability to identify the support of  $\mu$  by letting  $\varepsilon$  go to zero. However, we go beyond pointwise convergence and make an attempt to analyse whether weak convergence of a normalised version of  $q_{\varepsilon,n}(x)$  to  $\mu$  holds. Because  $q_{\varepsilon,n}$  is always larger than  $\varepsilon/N$ , it turns out that we have to consider the scaled function  $Nq_{\varepsilon,n}(x) - \varepsilon$  instead of  $q_{\varepsilon,n}$  for weak convergence.

<sup>67</sup>With more work, it is possible to derive bounds which maintain the interpolating property on the support. Nevertheless, we omit their presentation as these simple bounds suffice for our purposes.

**Example 3.3.6** (Single Dirac). If  $\mu = \delta_0$ , the regularised Christoffel function can be computed explicitly. The first singular function is simply  $N^{-1/2}D_{\text{rad},n/2}(x)$  while  $\sigma_1^{(n)} = N$ . Therefore, we have

$$Nq_{\varepsilon,n}(x) - \varepsilon = \frac{\varepsilon F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right)}{N - F_{\text{rad},n}(x) + \frac{\varepsilon}{N} F_{\text{rad},n}(x)} \quad (3.29)$$

for any  $x \in \mathbb{T}^d$ . In the following, we want to study whether  $\frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}} \rightarrow \mu$  as  $\varepsilon \rightarrow 0$  and thus estimate the distance of the two measures in the 1-Wasserstein metric.

- (i) Because there is only one possible transport map which pushes the mass from  $\frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}$  forward to  $\delta_0$ , we have

$$\begin{aligned} W_1 \left( \frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}, \delta_0 \right) &= \int_{\mathbb{T}^d} \frac{Nq_{\varepsilon,n}(x) - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}} \|x\|_{\mathbb{T}^d} dx \\ &= \|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}^{-1} \int_{[-\frac{1}{2}, \frac{1}{2}]^d} \frac{\varepsilon \|x\|_2 F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right) dx}{N - F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right)}. \end{aligned}$$

We need a lower bound for  $N - F_{\text{rad},n}(x)$  around the origin which is difficult to analyse directly. However, we observe that we have  $N - F_{\text{rad},n}(x) = N(1 - p_{1,n}(x))$  such that we can use the results from the previous section. Hence, we can fix some  $\tau > 0$  and have some constants  $c_{d,\tau}^{(9)}, c_{d,\tau}^{(11)} > 0$  with

$$N - F_{\text{rad},n}(x) \geq \begin{cases} c_{d,\tau}^{(9)} N n^2 \|x\|_2^2, & \text{if } \|x\|_2 \leq \frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi n} \\ c_{d,\tau}^{(11)} N, & \text{if } \|x\|_2 \geq \frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi n} \end{cases}$$

by Theorem 3.2.6 and Proposition 2.3.2. Then, using Lemma 1.3.7, we bound

$$\begin{aligned} &\int_{[-\frac{1}{2}, \frac{1}{2}]^d} \frac{\varepsilon \|x\|_2 F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right) dx}{N - F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right)} \\ &\leq \left(1 - \frac{\varepsilon}{N}\right) \left[ \int_0^{\frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi n}} \frac{c_d \varepsilon N r^d}{c_{d,\tau}^{(9)} (N - \varepsilon) n^2 r^2 + \varepsilon} dr + \int_{\frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi n}}^{\sqrt{d}/2} \frac{c_d \varepsilon N n^{-2} r^{d-2}}{c_{d,\tau}^{(11)} (N - \varepsilon) + \varepsilon} dr \right] \\ &\leq c_d \varepsilon \left(1 - \frac{\varepsilon}{N}\right) \begin{cases} c_\tau^{(1)} \left( \frac{\log\left(1 + \frac{c_\tau(N-\varepsilon)}{\varepsilon}\right)}{2n(N-\varepsilon)} + \frac{N}{c_{d,\tau}^{(11)}(N-\varepsilon)+\varepsilon} \frac{\log(c_\tau n)}{n^2} \right), & d = 1, \\ c_\tau^{(2)} \left( \frac{1}{(N-\varepsilon)n^2} + \frac{N}{c_{d,\tau}^{(11)}(N-\varepsilon)+\varepsilon} \frac{1}{n^2} \right), & d \geq 2, \end{cases} \\ &\leq c_d \varepsilon \left(1 - \frac{\varepsilon}{N}\right) \begin{cases} c_\tau^{(3)} \frac{\log\left(1 + \frac{c_\tau n}{\varepsilon}\right)}{n^2}, & d = 1, \\ c_\tau^{(4)} \frac{1}{n^2}, & d \geq 2, \end{cases} \end{aligned}$$

for some constants  $c_d, c_\tau, c_\tau^{(1)}, c_\tau^{(2)}, c_\tau^{(3)}, c_\tau^{(4)} > 0$ .

- (ii) On the other hand, with Remark 3.2.7 a lower bound on the total mass of  $Nq_{\varepsilon,n} - \varepsilon$  is given by

$$\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)} = \int_{[-\frac{1}{2}, \frac{1}{2}]^d} \frac{\varepsilon \left(1 - \frac{\varepsilon}{N}\right) F_{\text{rad},n}(x)}{N - \left(1 - \frac{\varepsilon}{N}\right) F_{\text{rad},n}(x)} dx$$

### 3 Trigonometric polynomials and rational functions

$$\begin{aligned}
&\geq \int_{B_{\frac{1}{\sqrt{2\pi dn}}}(0)} \frac{\varepsilon \left(1 - \frac{\varepsilon}{N}\right) N(1 - 2\pi^2 d^2 n^2 \|x\|_2^2)}{2\pi^2(N - \varepsilon)d^2 n^2 \|x\|_2^2 + \varepsilon} dx \\
&= c_d \int_0^1 \frac{\varepsilon \left(1 - \frac{\varepsilon}{N}\right) (1 - r^2)r^{d-1}}{(N - \varepsilon)r^2 + \varepsilon} dr \\
&\geq \frac{1}{4} c_d \int_0^{\frac{1}{2}} \frac{\varepsilon \left(1 - \frac{\varepsilon}{N}\right) r^{d-1}}{(N - \varepsilon)r^2 + \varepsilon} dr.
\end{aligned}$$

Interestingly, this integral has a very different behaviour in  $\varepsilon$  as  $\varepsilon \rightarrow 0$  depending on the dimension  $d$ . This is because the integrand has a pole in zero for  $\varepsilon \rightarrow 0$  if  $d \in \{1, 2\}$  whereas the singularity is removable if  $d \geq 3$ . In the latter case, the lower bound is proportional to  $\frac{\varepsilon}{N-\varepsilon} \int_0^{\frac{1}{2}} r^{d-3} - \frac{\varepsilon}{(N-\varepsilon)r^2+\varepsilon} dr$  and in general we end up with

$$\frac{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}{1 - \frac{\varepsilon}{N}} \geq \begin{cases} \frac{c_d}{4} \sqrt{\frac{\varepsilon}{N-\varepsilon}} \arctan\left(\frac{1}{2} \sqrt{\frac{N-\varepsilon}{\varepsilon}}\right), & d = 1, \\ \frac{c_d}{4} \frac{\varepsilon}{N-\varepsilon} \log\left(1 + \frac{N-\varepsilon}{4\varepsilon}\right), & d = 2, \\ \frac{c_d \varepsilon \left(\tilde{c}_d - \varepsilon^{1/2} (N-\varepsilon)^{-1/2} \arctan\left(\frac{\sqrt{N-\varepsilon}}{2\sqrt{\varepsilon}}\right)\right)}{4(N-\varepsilon)}, & d \geq 3, \end{cases} \quad (3.30)$$

for some  $\tilde{c}_d > 0$ .

(iii) Bringing (i) and (ii) together, we get

$$W_1\left(\frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}, \delta_0\right) \leq c_{d,\tau} \cdot \begin{cases} \frac{\sqrt{\varepsilon} \log\left(1 + \frac{c_\tau n}{\varepsilon}\right)}{\arctan\left(\frac{\sqrt{n-\varepsilon}}{2\sqrt{\varepsilon}}\right) n^{3/2}}, & d = 1, \\ \frac{1}{\log\left(1 + \frac{n^2-\varepsilon}{4\varepsilon}\right)}, & d = 2, \\ \frac{n^{d-2}}{\tilde{c}_d - \varepsilon^{1/2} (N-\varepsilon)^{-1/2} \arctan\left(\frac{(N-\varepsilon)^{1/2}}{2\sqrt{\varepsilon}}\right)}, & d \geq 3, \end{cases}$$

for some constant  $c_{d,\tau} > 0$ . This yields  $\frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}} \rightarrow \mu = \delta_0$  as  $\varepsilon \rightarrow 0$  for dimension  $d = \{1, 2\}$ .

(iv) On the contrary, we have by the same steps as in (ii) the lower bound

$$\begin{aligned}
\int_{[-\frac{1}{2}, \frac{1}{2}]^d} \frac{\varepsilon \|x\|_2 F_{\text{rad},n}(x)}{N - F_{\text{rad},n}(x) \left(1 - \frac{\varepsilon}{N}\right)} &\geq \frac{1}{4} c_d \int_0^{\frac{1}{2}} \frac{\varepsilon r^d}{(N - \varepsilon)r^2 + \varepsilon} dr \\
&\geq \frac{c_{d+1} \varepsilon \left(\tilde{c}_{d+1} - \varepsilon^{1/2} (N - \varepsilon)^{-1/2} \arctan\left(\frac{\sqrt{N-\varepsilon}}{2\sqrt{\varepsilon}}\right)\right)}{4(N - \varepsilon)}
\end{aligned}$$

in  $d \geq 3$ . Similarly, one can follow the lines of (i) in order to deduce that

$$\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)} \leq \varepsilon \left(1 - \frac{\varepsilon}{N}\right) c_\tau^{(4)} \frac{1}{n^2}$$

and this implies

$$W_1\left(\frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}, \delta_0\right) \geq cn^{2-d} \quad \text{for some constant } c > 0$$



such that the diverging rate in (iii) is sharp in  $\varepsilon$  giving  $W_1 \left( \frac{Nq_{\varepsilon,n} - \varepsilon}{\|Nq_{\varepsilon,n} - \varepsilon\|_{L^1(\mathbb{T}^d)}}, \mu \right) \rightarrow 0$  as  $\varepsilon \rightarrow 0$  for  $d \geq 3$ . On the other hand, the question of convergence for  $n \rightarrow \infty$  needs a more careful analysis because the upper and lower bounds show a different convergence behaviour in  $n$ .

As we have seen in Example 3.3.6 that we can control the situation for one node with sufficient accuracy, it might be natural to consider a relation between  $Nq_{\varepsilon,n} - \varepsilon$  for multiple nodes and the sum of the individual Christoffel functions. In fact, we develop the following bound through the variational formulation of the Christoffel function.

**Lemma 3.3.7** (Lower bound on  $q_{\varepsilon,n}$ ). *Let  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  with  $\alpha_t > 0$ ,  $\sum_t \alpha_t = 1$  and  $\text{sep } Y = \frac{2\sqrt{1+\tau}jd/2-1,1}{\pi(n-1)}$  for some  $\tau > 0$  and  $\frac{\varepsilon}{N \min_t \alpha_t} < 1 + \|\mathbf{I}_r - (N^{-1} \mathbf{A}_n^* \mathbf{A}_n)^{-1}\|_2$ . Then, for every  $x \in \mathbb{T}^d$  the lower bound*

$$Nq_{\varepsilon,n}(x) - \varepsilon \geq \sum_{t \in Y} \alpha_t \frac{\varepsilon \left( 1 - \frac{\varepsilon}{N} \left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right]^{-1} \right) F_{\text{rad},n}(x-t)}{N - F_{\text{rad},n}(x-t) + \frac{\varepsilon}{N} \left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right]^{-1} F_{\text{rad},n}(x-t)}$$

is valid.

*Proof.* We have to make the variational formulation (3.28) more explicit. Therefore, let  $\mathbf{U}_{1,n} \in \mathbb{C}^{N \times |Y|}$  be the matrix whose columns are the first  $|Y|$  singular vectors of  $\mathbf{T}_n$  and note that  $\mathbf{I}_N = \mathbf{U}_n \mathbf{U}_n^* = \mathbf{U}_{1,n} \mathbf{U}_{1,n}^* + \mathbf{U}_{0,n} \mathbf{U}_{0,n}^*$ . By  $\text{img } \mathbf{A}_n = \text{img } \mathbf{U}_{1,n}$ , there is a regular matrix  $\mathbf{B}$  such that  $\mathbf{U}_{1,n} = \mathbf{A}_n \mathbf{B}$  and thus  $\mathbf{I}_{|Y|} = \mathbf{U}_{1,n}^* \mathbf{U}_{1,n} = \mathbf{B}^* \mathbf{A}_n^* \mathbf{A}_n \mathbf{B}$  or  $(\mathbf{B} \mathbf{B}^*)^{-1} = \mathbf{A}_n^* \mathbf{A}_n$  respectively. Hence, we can rewrite the objective function in (3.28) as

$$\begin{aligned} \mathbf{p}^* (\mathbf{T}_n + \varepsilon \mathbf{U}_{0,n} \mathbf{U}_{0,n}^*) \mathbf{p} &= \mathbf{p}^* \left( \mathbf{T}_n - \varepsilon \mathbf{A}_n (\mathbf{A}_n^* \mathbf{A}_n)^{-1} \mathbf{A}_n^* + \varepsilon \mathbf{I}_N \right) \mathbf{p} \\ &= \mathbf{p}^* \mathbf{A}_n \left( \mathbf{W} - \frac{\varepsilon}{N} \mathbf{I}_r + \varepsilon \left( N^{-1} \mathbf{I}_r - (\mathbf{A}_n^* \mathbf{A}_n)^{-1} \right) \right) \mathbf{A}_n^* \mathbf{p} + \varepsilon \|\mathbf{p}\|_2^2 \\ &\geq \mathbf{p}^* \mathbf{A}_n \left( \mathbf{W} - \frac{\varepsilon \left( 1 + \|\mathbf{I}_r - (N^{-1} \mathbf{A}_n^* \mathbf{A}_n)^{-1}\|_2 \right)}{N} \mathbf{I}_r \right) \mathbf{A}_n^* \mathbf{p} + \varepsilon \|\mathbf{p}\|_2^2 \\ &\geq \sum_t \alpha_t \left( \left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right] |e_t^* \mathbf{A}_n^* \mathbf{p}|^2 + \varepsilon \|\mathbf{p}\|_2^2 \right) \end{aligned}$$

where we used the normalisation of the weights and abbreviate the term  $c_{n,d,\text{sep } Y} := \|\mathbf{I}_r - (N^{-1} \mathbf{A}_n^* \mathbf{A}_n)^{-1}\|_2 \in \mathcal{O}(n^{-1})$ , see the proof of Lemma 3.2.8. Then, the variational formulation (3.28) and the representation for a single Dirac (3.29) allow to bound

$$\begin{aligned} q_{\varepsilon,n}(x) &\geq \min \left\{ \sum_t \alpha_t \left( \left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right] |e_t^* \mathbf{A}_n^* \mathbf{p}|^2 + \varepsilon \|\mathbf{p}\|_2^2 \right) : p \in \mathcal{P}^{n/2,d,2}, p(x) = 1 \right\} \\ &\geq \sum_{t \in Y} \alpha_t \min \left\{ \left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right] |e_t^* \mathbf{A}_n^* \mathbf{p}|^2 + \varepsilon \|\mathbf{p}\|_2^2 : p \in \mathcal{P}^{n/2,d,2}, p(x) = 1 \right\} \\ &= \sum_{t \in Y} \alpha_t \frac{1}{\left[ 1 - \frac{\varepsilon(1+c_{n,d,\text{sep } Y})}{\alpha_t N} \right]^{-1} N^{-1} F_{\text{rad},n}(x-t) + \varepsilon^{-1} (N - F_{\text{rad},n}(x-t))}. \end{aligned}$$

such that the proposed bound follows directly.  $\square$

### 3 Trigonometric polynomials and rational functions

Even though this Lemma gives an idea how it might be possible to link the regularised Christoffel function of a discrete measure to the sum of the individual Christoffel functions, it is not clear whether  $\frac{Nq_{\varepsilon,n}-\varepsilon}{\|Nq_{\varepsilon,n}-\varepsilon\|_{L^1}} \rightharpoonup \mu$  as  $\varepsilon \rightarrow 0$  holds for general  $\mu \in \mathcal{M}(q)$ . By Example 3.3.6, we know that this can be only true for the univariate or bivariate case. Moreover, the bounds from this section help us to choose  $\varepsilon$  depending on the noise level if we deal with inexact moments.

## 3.4 Recovery from noisy data

Assume we are given noisy data  $\hat{\mu} = \hat{\mu}_0 + \hat{\rho}$  originating from some perturbation of the ground truth  $\mu_0 = \sum_{t \in Y} \alpha_t \delta_t \in \mathcal{M}(q) \cap \mathcal{M}_{+,1}(\mathbb{T}^d)$  such that we have access to the moment matrix

$$\mathbf{T}_n = (\hat{\mu}(k-l))_{k,l \in [n]} = \mathbf{T}_{0,n} + \mathbf{T}_{\varrho,n},$$

where  $\mathbf{T}_{0,n}$  is the moment matrix corresponding to the ground truth measure  $\mu_0$  and  $\mathbf{T}_{\varrho,n} = (\hat{\rho}(k-l))_{k,l \in [n]}$  contains the noise with  $|\hat{\rho}(k)| \leq \varrho$  for some  $\varrho > 0$ . In this section we want to use the rational approximation of the previous Section 3.3 and develop a theory how the regularisation parameter  $\varepsilon > 0$  should be chosen in order to obtain an optimal approximation to the unknown support of the ground truth measure  $\mu_0$ . To the best of our knowledge, such a recovery of a discrete measure from noisy data using Christoffel functions has not been analysed before.<sup>68</sup>

When we compute the SVD  $\mathbf{T}_n = \tilde{\mathbf{V}}_n \tilde{\Sigma}_n \tilde{\mathbf{V}}_n^*$ , we cannot expect that this matrix satisfies  $\text{rank } \mathbf{T}_n = |Y|$ . Instead, the matrix is full rank and we need to truncate the smallest singular values which are dominated by the noise. Therefore, it is natural to take the smallest  $\tilde{r} \in \mathbb{N}$  such that the best  $\tilde{r}$ -term approximation  $\mathbf{T}_{\tilde{r},n} = \tilde{\mathbf{V}}_n \text{diag}(\tilde{\sigma}_1^{(n)}, \dots, \tilde{\sigma}_{\tilde{r}}^{(n)}, 0, \dots, 0) \tilde{\mathbf{V}}_n^*$  satisfies

$$\|\mathbf{T}_{\tilde{r},n} - \mathbf{T}_n\|_2 = \tilde{\sigma}_{\tilde{r}+1}^{(n)} \leq \|\mathbf{T}_{\varrho,n}\|_2 \leq N\varrho \quad (3.31)$$

and hope that  $\mathbf{T}_{\tilde{r},n}$  is close to the ground truth  $\mathbf{T}_{0,n}$ .<sup>69,70</sup> In particular, the idea is that  $\tilde{r} = |Y|$  if the noise level  $\varrho$  is sufficiently small and that the first singular values and vectors are close to the original ones. However, this approach is based on a very accurate knowledge of the size of the noise. Instead, we generalise the concept of rational approximation from Section 3.3 to the noisy case and study the regularised Christoffel function

$$\tilde{q}_{\varepsilon,n}(x) = \frac{1}{\sum_{j=1}^N g_{\varepsilon}(\tilde{\sigma}_j^{(n)}) |\tilde{u}_j^{(n)}(x)|^2} = \frac{1}{\sum_{j:\tilde{\sigma}_j^{(n)} > \varepsilon} \frac{|\tilde{u}_j^{(n)}(x)|^2}{\tilde{\sigma}_j^{(n)}} + \sum_{j:\tilde{\sigma}_j^{(n)} \leq \varepsilon} \frac{|\tilde{u}_j^{(n)}(x)|^2}{\varepsilon}} \quad (3.32)$$

with regularisation parameter  $\varepsilon > 0$ . This comes with the benefit that this expression can be computed without any knowledge of the noise level. Nevertheless, the parameter  $\varepsilon$  must

<sup>68</sup>In [111, Sec. 3.3], an underlying measure which is absolutely continuous with respect to Lebesgue measure is approximated by a regularised Christoffel function.

<sup>69</sup>In applications, it is more natural to assume knowledge about the noise level  $\varrho$  than on  $\|\mathbf{T}_{\varrho,n}\|_2$ . Hence, the numerical rank  $\tilde{r}$  will be chosen such that  $\tilde{\sigma}_{\tilde{r}+1}^{(n)} \leq N\varrho$  in these situations. However, this does not change the situation if  $\tilde{\sigma}_{\tilde{r}}^{(n)} \geq \sigma_{\tilde{r}}^{(n)} - \|\mathbf{T}_{\varrho,n}\|_2 \geq N(\alpha_{\min} - \varrho) - \mathcal{O}(n^{d-1}) > N\varrho$  i.e. if  $\alpha_{\min} > 2\varrho$  and  $N$  sufficiently large.

<sup>70</sup>Note that  $\mathbf{T}_{\tilde{r},n}$  is in general not a Toeplitz matrix. Rank-1 best approximations of matrices by Toeplitz (or equivalently Hankel) matrices has been studied in [82] but the case of a larger rank is difficult, see [81]. Therefore, we use the unstructured approximation of  $\mathbf{T}_n$  by  $\mathbf{T}_{\tilde{r},n}$ .

be tuned depending on the noise level and it is the aim of this section to develop an idea how one should choose  $\varepsilon$  in order to obtain a Christoffel function  $\tilde{q}_{\varepsilon,n}$  which peaks around the support of  $\mu_0$ . At first, we study under which condition on the noise the truncation via (3.31) recovers the correct number of parameters.

**Lemma 3.4.1.** *Let the smallest weight  $\alpha_{\min} = \min_t \alpha_t$  satisfy*

$$\alpha_{\min} > \frac{2\|\mathbf{T}_{\varrho,n}\|_2}{N} + \frac{c_{d,\tau}^{(10)}}{n} \cdot \frac{\alpha_{\max}|Y|}{\text{sep } Y}$$

where  $c_{d,\tau}^{(10)}$  is the constant from Lemma 3.2.8 and  $\tau > 0$  such that the separation fulfils  $\text{sep } Y \cdot (n-1) = \frac{2\sqrt{1+\tau}j_{d/2,1}}{\pi}$ . Then, we have  $\tilde{r} = |Y|$ .

*Proof.* By Weyl's bound on singular values, see Lemma 1.1.9, we can bound

$$\tilde{\sigma}_{|Y|}^{(n)} \geq \sigma_{|Y|}^{(n)} - \|\mathbf{T}_{\varrho,n}\|_2 \geq N \left( \alpha_{\min} - \frac{c_{d,\tau}^{(10)}}{n} \cdot \frac{\alpha_{\max}|Y|}{\text{sep } Y} \right) - \|\mathbf{T}_{\varrho,n}\|_2 > \|\mathbf{T}_{\varrho,n}\|_2$$

where we use Lemma 3.2.8 in the second inequality. On the other hand, Weyl's inequality also gives  $\tilde{\sigma}_{|Y|+1}^{(n)} \leq \|\mathbf{T}_{\varrho,n}\|_2$  and thus  $\tilde{r} = |Y|$ .  $\square$

The assumption of the previous lemma is natural and can be seen as a lower bound on the *signal-to-noise ratio (SNR)* which is a classical assumption in signal analysis, e.g. see [103]. Next, we analyse how the signal polynomial is perturbed by the noise.

**Lemma 3.4.2.** *Under the conditions of Lemma 3.4.1, we have that the noisy version  $\tilde{p}_{1,n} = \frac{1}{N} \sum_{j=1}^{|Y|} |\tilde{u}_j^{(n)}(x)|^2$  of the signal polynomial satisfies*

$$|p_{1,n}(x) - \tilde{p}_{1,n}(x)| \leq \frac{2\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \tilde{p}_{1,n}(x) \max_j \frac{|\tilde{\sigma}_j^{(n)} - \sigma_j^{(n)}|}{\sigma_j^{(n)}} \leq \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}}.$$

*Proof.* As  $\text{rank } \mathbf{T}_n = |Y|$  by Lemma 3.4.1, we have

$$\begin{aligned} |p_{1,n}(x) - \tilde{p}_{1,n}(x)| &\leq \left| \sum_{j=1}^{|Y|} \frac{\sigma_j^{(n)} |u_j^{(n)}(x)|^2 - \tilde{\sigma}_j^{(n)} |\tilde{u}_j^{(n)}(x)|^2}{N\sigma_j^{(n)}} \right| + \left| \sum_{j=1}^{|Y|} \frac{\tilde{\sigma}_j^{(n)} - \sigma_j^{(n)}}{N\sigma_j^{(n)}} |\tilde{u}_j^{(n)}(x)|^2 \right| \\ &\leq \frac{1}{\sigma_{|Y|}^{(n)}} \|\mathbf{T}_{\tilde{r},n} - \mathbf{T}_{0,n}\|_2 + \tilde{p}_{1,n}(x) \max_j \frac{|\tilde{\sigma}_j^{(n)} - \sigma_j^{(n)}|}{\sigma_j^{(n)}} \\ &\leq \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \end{aligned}$$

where the last inequality follows from Weyl's inequality for the perturbation of singular values, cf. Lemma 1.1.9.  $\square$

Additionally, we remark that by definition the perturbed noise-polynomial  $\tilde{p}_{0,n}$  admits  $\tilde{p}_{0,n}(x) = 1 - \tilde{p}_{1,n}$  such that the previous Lemma 3.4.2 holds also for  $\tilde{p}_{0,n}$ . Analogously, we can derive a result on the difference between  $\tilde{q}_{\varepsilon,n}$  and the exact Christoffel function  $q_{\varepsilon,n}$  corresponding to the ground truth  $\mu_0$ .

### 3 Trigonometric polynomials and rational functions

**Lemma 3.4.3.** *If  $\varepsilon < \min(\sigma_{|Y|}^{(n)}, \tilde{\sigma}_{|Y|}^{(n)})$  and with the conditions on the noise level from Lemma 3.4.1, we have that the noisy version  $\tilde{q}_{\varepsilon,n}$  of the Christoffel function and  $q_{\varepsilon,n}$  being the Christoffel function associated to  $\mu$  satisfy*

$$\frac{|q_{\varepsilon,n}(x) - \tilde{q}_{\varepsilon,n}(x)|}{\left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right) q_{\varepsilon,n}(x)} \leq \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \left( \frac{N}{\varepsilon} + \frac{1 + \sqrt{5}}{2\tilde{\sigma}_{|Y|}^{(n)}} N \right) + \max(0, \varepsilon^{-1} - \|\mathbf{T}_{\varrho,n}\|_2^{-1})N$$

for any  $\varepsilon > 0$  and  $x \in \mathbb{T}^d$ .

*Proof.* Using  $\varepsilon < \min(\sigma_{|Y|}^{(n)}, \tilde{\sigma}_{|Y|}^{(n)})$ , we estimate

$$\begin{aligned} & \frac{|q_{\varepsilon,n}(x) - \tilde{q}_{\varepsilon,n}(x)|}{q_{\varepsilon,n}(x)} \\ &= \left| \frac{\sum_{j=1}^N g_{\varepsilon}(\sigma_j^{(n)}) |u_j^{(n)}(x)|^2 - g_{\varepsilon}(\tilde{\sigma}_j^{(n)}) |\tilde{u}_j^{(n)}(x)|^2}{\sum_{j=1}^N g_{\varepsilon}(\tilde{\sigma}_j^{(n)}) |\tilde{u}_j^{(n)}(x)|^2} \right| \\ &\leq \frac{\left| \sum_{j=1}^{|Y|} \left( \frac{1}{\sigma_j^{(n)}} - \varepsilon^{-1} \right) |u_j^{(n)}(x)|^2 - \sum_{j: \tilde{\sigma}_j^{(n)} > \varepsilon} \left( \frac{1}{\tilde{\sigma}_j^{(n)}} - \varepsilon^{-1} \right) |\tilde{u}_j^{(n)}(x)|^2 \right|}{\frac{1}{\tilde{\sigma}_1^{(n)}} N} \\ &\leq \frac{\left| \sum_{j=1}^{|Y|} \frac{|u_j^{(n)}(x)|^2}{\sigma_j^{(n)}} - \frac{|u_j^{(n)}(x)|^2}{\varepsilon} - \frac{|\tilde{u}_j^{(n)}(x)|^2}{\tilde{\sigma}_j^{(n)}} + \frac{|\tilde{u}_j^{(n)}(x)|^2}{\varepsilon} \right| + \sum_{j: \varepsilon < \tilde{\sigma}_j^{(n)} < \|\mathbf{T}_{\varrho,n}\|_2} \frac{|\tilde{u}_j^{(n)}(x)|^2}{\varepsilon} - \frac{|\tilde{u}_j^{(n)}(x)|^2}{\tilde{\sigma}_j^{(n)}}}{\left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right)^{-1}} \\ &\leq \frac{\frac{N}{\varepsilon} |p_{1,n}(x) - \tilde{p}_{1,n}(x)| + \frac{1+\sqrt{5}}{2\sigma_{|Y|}^{(n)} \tilde{\sigma}_{|Y|}^{(n)}} N \|\mathbf{T}_{0,n} - \mathbf{T}_{|Y|,n}\|_2 + \max(0, \varepsilon^{-1} - \|\mathbf{T}_{\varrho,n}\|_2^{-1})N}{\left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right)^{-1}} \\ &\leq \left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right) \left[ \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \left( \frac{N}{\varepsilon} + \frac{1 + \sqrt{5}}{2\tilde{\sigma}_{|Y|}^{(n)}} N \right) + \max(0, \varepsilon^{-1} - \|\mathbf{T}_{\varrho,n}\|_2^{-1})N \right] \end{aligned}$$

and remark that the first inequality uses the assumption on  $\varepsilon$ , the second bounds the singular value  $\tilde{\sigma}_1^{(n)} \leq N + \|\mathbf{T}_{\varrho,n}\|_2$ , the third applies Theorem 1.1.10 and the fourth takes Lemma 3.4.2 into account.  $\square$

We now come back to the original problem to derive a parameter choice rule for  $\varepsilon$ . At first, it might seem natural to choose  $\varepsilon > 0$  such that the upper bound on the perturbed Christoffel function  $\tilde{q}_{\varepsilon,n}$  given by

$$\tilde{q}_{\varepsilon,n}(x) \leq q_{\varepsilon,n}(x) + |q_{\varepsilon,n}(x) - \tilde{q}_{\varepsilon,n}(x)| \quad (3.33)$$

becomes as small as possible for  $x$  outside of the support of  $\mu_0$ . Using Lemma 3.4.3, this upper bound on the perturbed Christoffel function becomes minimal for  $\varepsilon = 0$  and the optimal value of the bound agrees up to a constant with the bound for the signal polynomial presented in Lemma 3.4.2. Therefore, such a choice does not lead to a better result than the signal polynomial with the additional drawback that  $\varepsilon \rightarrow 0$  promotes singularity of the Christoffel function which is not beneficial for its visual representation.

The most natural approach might be to study the perturbation of the maxima of  $\tilde{q}_{\varepsilon,n}$  but this appears to be a complicated problem as it demands more than just an  $L^\infty$  bound

on the function. Instead, we propose to choose  $\varepsilon$  such that the fraction of the mass outside of the support in relation to the total mass of the Christoffel function becomes small. More precisely, we demand that for  $\delta = \frac{\text{sep}_Y}{2}$  and  $y \in \mathbb{T}^d$  such that  $B_\delta(y) \cap Y = \emptyset$  the normalised fraction

$$\frac{\delta^{-d} \int_{B_\delta(y)} N\tilde{q}_{\varepsilon,n}(x) - \varepsilon dx}{\int_{\mathbb{T}^d} N\tilde{q}_{\varepsilon,n}(x) - \varepsilon dx} \quad (3.34)$$

becomes as small as possible.<sup>71</sup> Here, we use again the scaled version  $N\tilde{q}_{\varepsilon,n} - \varepsilon$  as the analysis of weak convergence has already shown the benefit of this approach which is that the integral over  $N\tilde{q}_{\varepsilon,n} - \varepsilon$  is asymptotically different from  $\varepsilon/N$ . Under the assumption  $3\|\mathbf{T}_{\varrho,n}\|_2 < \frac{1}{2}\sigma_{|Y|}^{(n)}$  we introduce some  $\tilde{\delta}$  to be chosen such that

$$\tilde{\delta}^2 \pi^2 d^2 n^2 = \frac{1}{4} - \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{2\sigma_{|Y|}^{(n)}}.$$

Additionally, we assume the conditions of Lemma 3.4.3. Then, we bound the denominator of (3.34) for some constant  $c_d > 0$  as

$$\begin{aligned} \int_{\mathbb{T}^d} N\tilde{q}_{\varepsilon,n}(x) - \varepsilon dx &= \int_{\mathbb{T}^d} \frac{\varepsilon \sum_{j:\tilde{\sigma}_j^{(n)} > \varepsilon} \left(1 - \frac{\varepsilon}{\tilde{\sigma}_j^{(n)}}\right) |\tilde{u}_j^{(n)}(x)|^2}{N - \sum_{j:\tilde{\sigma}_j^{(n)} > \varepsilon} \left(1 - \frac{\varepsilon}{\tilde{\sigma}_j^{(n)}}\right) |\tilde{u}_j^{(n)}(x)|^2} dx \\ &\geq \int_{\mathbb{T}^d} \frac{\varepsilon \sum_{j=1}^{|Y|} \left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) |\tilde{u}_j^{(n)}(x)|^2}{N - \sum_{j=1}^{|Y|} \left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) |\tilde{u}_j^{(n)}(x)|^2} dx \\ &\geq \left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) \sum_{t \in Y} \int_{B_{\tilde{\delta}}(t)} \frac{\varepsilon \tilde{p}_{1,n}(x)}{1 - \tilde{p}_{1,n}(x) + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}} dx \\ &\geq \left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) \sum_{t \in Y} \int_{B_{\tilde{\delta}}(t)} \frac{\varepsilon \left(1 - 2\pi^2 d^2 n^2 \|x - t\|_2^2 - \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}}\right)}{2\pi^2 d^2 n^2 \|x - t\|_2^2 + \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}} dx \\ &\geq \frac{c_d |Y|}{2} \int_0^{\tilde{\delta}} \frac{\left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) \varepsilon r^{d-1}}{(2\pi^2 d^2 n^2) r^2 + \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}} dr \end{aligned}$$

<sup>71</sup>Similarly as in (3.33), we decompose (3.34) into an approximation error going to zero as  $\varepsilon \rightarrow 0$  and a regularisation error term which goes to infinity as  $\varepsilon \rightarrow 0$ . The latter is a typical phenomenon for the parameter choice of regularised inverse problems where one needs to balance between the two effects.

### 3 Trigonometric polynomials and rational functions

$$= \left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) \frac{\varepsilon c_d |Y|}{2} \begin{cases} \frac{\arctan \left( \sqrt{\frac{\frac{1}{2} - \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}}}{\frac{3\|\mathbf{T}_{\varrho,n}\|_2 + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}}}}{\sqrt{2\pi}dn \sqrt{\frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}}} \right)}{\log \left( 1 + \frac{\sqrt{\frac{1}{2} - \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}}}}{\frac{3\|\mathbf{T}_{\varrho,n}\|_2 + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}}} \right)}, & d = 1, \\ \frac{\log \left( 1 + \frac{\sqrt{\frac{1}{2} - \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}}}}{\frac{3\|\mathbf{T}_{\varrho,n}\|_2 + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}}} \right)}{4\pi^2 d^2 n^2}, & d = 2, \end{cases}$$

using the SVD representation from Lemma 3.3.5, whereas we can apply Lemma 3.3.5 and Lemma 3.4.3 in order to control the numerator of (3.34) by

$$\begin{aligned} & \int_{B_\delta(y)} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx \\ & \leq \left\{ 1 + \left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right) \left[ \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \left(\frac{N}{\varepsilon} + \frac{1+\sqrt{5}}{2\tilde{\sigma}_{|Y|}^{(n)}}N\right) + \max\left(0, \frac{N}{\varepsilon} - \frac{N}{\|\mathbf{T}_{\varrho,n}\|_2}\right) \right] \right\} \int_{B_\delta(y)} N q_{\varepsilon,n}(x) dx \\ & \leq \int_{B_\delta(y)} \frac{\left\{ 1 + \left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right) \left[ \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \left(\frac{N}{\varepsilon} + \frac{1+\sqrt{5}}{2\tilde{\sigma}_{|Y|}^{(n)}}N\right) + \max\left(0, \frac{N}{\varepsilon} - \frac{N}{\|\mathbf{T}_{\varrho,n}\|_2}\right) \right] \right\} \varepsilon}{c_{d,\tau}^{(9)} n^2 \text{dist}(x, Y)^2 + \frac{\varepsilon}{\sigma_1^{(n)}}} dx \\ & \leq c'_d \cdot \frac{\left\{ 1 + \left[ \frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} \left(\frac{1}{\varepsilon} + \frac{1+\sqrt{5}}{2\tilde{\sigma}_{|Y|}^{(n)}}\right) + \max\left(0, \frac{1}{\varepsilon} - \frac{1}{\|\mathbf{T}_{\varrho,n}\|_2}\right) \right] \left(1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{N}\right) \right\} \varepsilon}{n^2 \text{dist}(B_\delta(y), Y)^2} \end{aligned}$$

for some constant  $c'_d > 0$ . To sum up, we end up with the following bound.

**Proposition 3.4.4** (Optimal choice of  $\varepsilon$ ). *Let the conditions of Lemma 3.4.1 be fulfilled and additionally assume  $3\|\mathbf{T}_{\varrho,n}\|_2 < \frac{1}{4}\sigma_{|Y|}^{(n)}$  as well as  $\varepsilon \leq \min\left(\tilde{\sigma}_{|Y|}^{(n)}, \sigma_{|Y|}^{(n)}\right)$ .<sup>72</sup> Then, we can bound (3.34) from above by a dimension dependent constant times*

$$h(\varepsilon) := \frac{1 + \frac{\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}\varepsilon} + \max\left(0, \frac{1}{\varepsilon} - \frac{1}{\|\mathbf{T}_{\varrho,n}\|_2}\right)}{\left(1 - \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right) n \text{dist}(B_\delta(y), Y)^2 |Y|} \begin{cases} \sqrt{\frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}} & , d = 1, \\ \frac{n}{\log\left(1 + \left(\frac{3\|\mathbf{T}_{\varrho,n}\|_2}{\sigma_{|Y|}^{(n)}} + \frac{\varepsilon}{\tilde{\sigma}_{|Y|}^{(n)}}\right)^{-1}\right)} & , d = 2. \end{cases} \quad (3.35)$$

If we choose some rate  $\varepsilon = \|\mathbf{T}_{\varrho,n}\|_2^a$ ,  $a > 0$  for the regularisation parameter, the resulting rate for the bound  $h$  as  $\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0$  is

$$\sup \left\{ b : \lim_{\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0} \frac{h(\|\mathbf{T}_{\varrho,n}\|_2^a)}{\|\mathbf{T}_{\varrho,n}\|_2^b} < \infty \right\} = \begin{cases} \frac{1}{2} & , a = 1, \\ \frac{1}{2}a & , a \in (0, 1), \\ \frac{1}{2} - a & , a > 1, \end{cases}$$

for dimension  $d = 1$  and

$$\sup \left\{ b : \lim_{\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0} \frac{h(\|\mathbf{T}_{\varrho,n}\|_2^a)}{\|\mathbf{T}_{\varrho,n}\|_2^b} < \infty \right\} = \begin{cases} 0 & , a \in (0, 1], \\ -a + 1 & , a > 1, \end{cases}$$

<sup>72</sup>Note that it is reasonable to take  $\varepsilon$  smaller than the smallest meaningful singular value of  $\mathbf{T}_n$ .

for  $d = 2$ . Therefore, the order-optimal choice is  $\varepsilon^* = \|\mathbf{T}_{\varrho,n}\|_2$  leading to the optimal value

$$h(\varepsilon^*) = \|\mathbf{T}_{\varrho,n}\|_2^{1/2} \cdot \mathcal{O}\left(n^{-3/2}\right)$$

for  $d = 1$ , whereas any exponent  $a \in (0, 1]$  is optimal in the bivariate setting. Choosing  $a = 1$  for  $d = 2$  gives the pointwise bound

$$\frac{\delta^{-d} \int_{B_\delta(y)} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx}{\int_{\mathbb{T}^2} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx} \leq \frac{c_2}{\log(1 + n^2 \|\mathbf{T}_{\varrho,n}\|_2^{-1})}$$

for some constant  $c_2$  depending on  $Y, y, \alpha_{\min}$  and  $\alpha_{\max}$ .

*Proof.* The previously derived bounds on (3.34) are the first ingredient for (3.35) which then follows from the (reasonable) assumptions on  $\|\mathbf{T}_{\varrho,n}\|_2$  and  $\varepsilon$  in relation to  $\tilde{\sigma}_Y^{(n)}$ . This includes the observation that the argument of arctan is larger than  $\frac{1}{2}\sqrt{\frac{4}{5}}$  such that we can discard the contribution of arctan as a constant because of the monotony of arcus tangens. The influence of the exponent  $a$  on the convergence rate of  $h(\|\mathbf{T}_{\varrho,n}\|_2^a)$  as  $\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0$  is straightforward. Finally, the rates for  $h(\|\mathbf{T}_{\varrho,n}\|_2)$  in terms of  $\|\mathbf{T}_{\varrho,n}\|_2$  and  $n$  can be directly concluded from Lemma 3.2.8 where we have seen that  $\sigma_j^{(n)} \in \mathcal{O}(N) = \mathcal{O}(n^d)$ .  $\square$

We remark that it appears difficult to find the exact minimum of  $h$  as in (3.35) which is the reason why we circumvented the problem by simplifying it to the question of the optimal order as  $\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0$ . As a consequence of this approach, it might be a good choice in applications to take  $\varepsilon$  proportional to  $\|\mathbf{T}_{\varrho,n}\|_2$  for both  $d = 1$  and  $d = 2$ . With this asymptotic parameter choice, Proposition 3.4.4 gives a convergence rate for

$$\frac{\delta^{-d} \int_{B_\delta(y)} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx}{\int_{\mathbb{T}^2} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx} \rightarrow 0$$

if  $\|\mathbf{T}_{\varrho,n}\|_2 \rightarrow 0$  or  $n \rightarrow \infty$ .

**Example 3.4.5.** In order to validate that this parameter choice rule for  $\varepsilon$  indeed minimises our objective (3.34) at least in a simple example, we perturb the data on which Figure 3.6 is based.<sup>73</sup> More precisely, we vary the size of the noise  $\|\mathbf{T}_{\varrho,n}\|_2$  and set up the Christoffel function  $\tilde{q}_{\varepsilon,n}$  for fixed  $n = 24$  and various  $\varepsilon \in [0, 6]$ . Here, the perturbed moments  $\hat{\rho}(k)$  are obtained by independent and identically distributed samples of a normally distributed random variable and after storing them as entries of the Toeplitz matrix  $\mathbf{T}_{\varrho,n}$  subsequent normalisation to a fixed value of  $\|\mathbf{T}_{\varrho,n}\|_2$  is performed. For each of these configurations, we approximately compute

$$\frac{\int_{[0,0.08] \cup [0.45,0.58] \cup [0.9,1]} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx}{\int_{\mathbb{T}^2} N \tilde{q}_{\varepsilon,n}(x) - \varepsilon dx} \quad (3.36)$$

in order to estimate how small the Christoffel function  $\tilde{q}_{\varepsilon,n}$  is outside of the support of the actual measure  $\mu_0$  with  $\text{supp } \mu_0 = \{\frac{1}{6}, \frac{1}{3}, \frac{3}{4}\}$ . Additionally, the same analysis is performed for fixed perturbation  $\|\mathbf{T}_{\varrho,n}\|_2 = \varepsilon = 1.5$  and  $n \in \{5, 6, \dots, 150\}$  in order to check the rate of  $\mathcal{O}(n^{-3/2})$  presented for the contribution away from the support in Proposition 3.4.4. In

<sup>73</sup>Apart from the code described in Subsection 3.2.3, all code which is relevant for this work can be found on GitHub, see <https://github.com/MHockmann/Dissertation.git>.

### 3 Trigonometric polynomials and rational functions

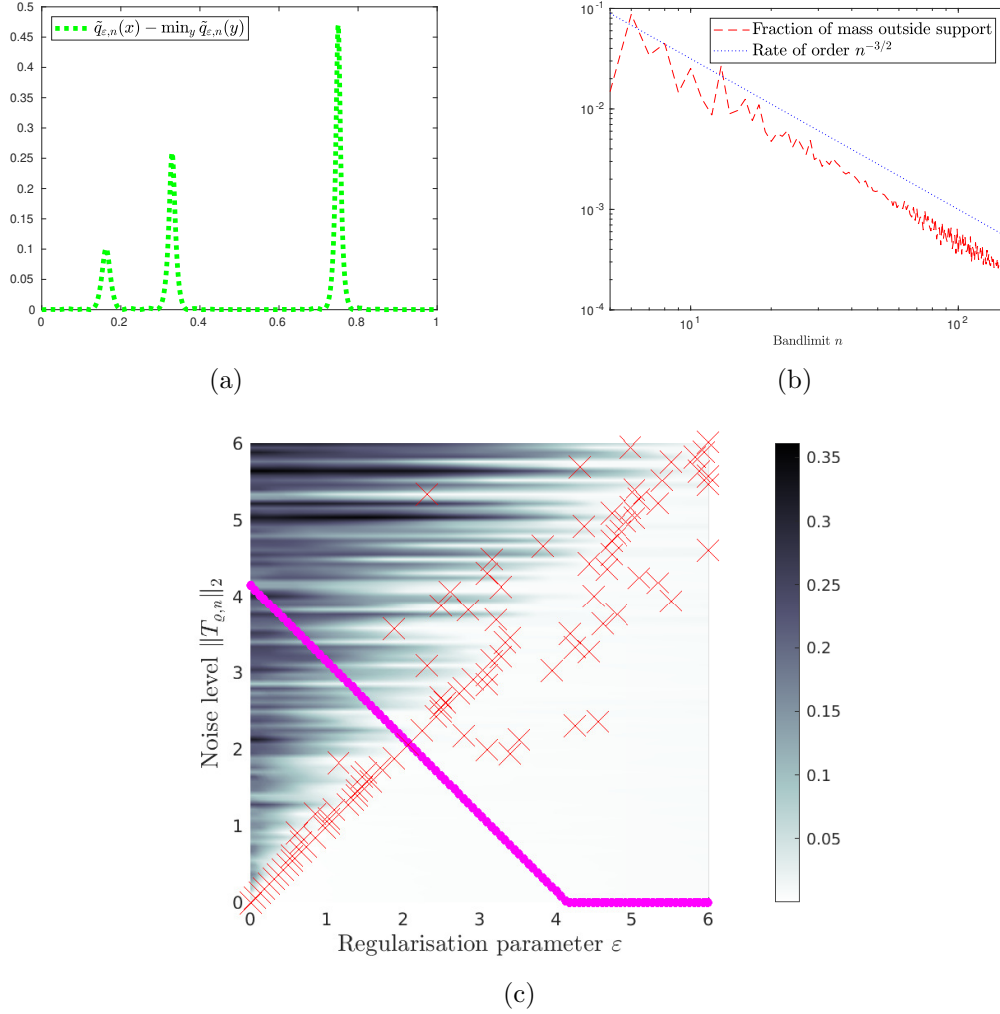


Figure 3.7: Parameter choice for perturbed Christoffel function. The exact  $\mathbf{q}_{\varepsilon,n}$  to the ground truth  $\mu_0$  was considered in Figure 3.6. In (a) we display  $\tilde{q}_{\varepsilon,n}(x) - \min_y \tilde{q}_{\varepsilon,n}(y)$  for  $n = 24$  and  $\varepsilon = \|\mathbf{T}_{\varrho,n}\|_2 = 1.5$  while (b) compares (3.36) to  $n^{-3/2}$  for  $\|\mathbf{T}_{\varrho,n}\|_2 = \varepsilon = 1.5$  and  $n \in \{5, 6, \dots, 150\}$ . The optimal choice of  $\varepsilon$  balancing effects if  $\varepsilon \rightarrow 0$  and  $\varepsilon \rightarrow \sigma_3^{(n)}$  is shown in (c) where for fixed  $n = 24$  and various  $\varepsilon$ ,  $\|\mathbf{T}_{\varrho,n}\|_2 \in [0, 6]$  we illustrate the value of (3.36). Moreover, we highlight for each noise level the minimiser  $\varepsilon^*$  of (3.36) (red crosses) and add the barrier  $\|\mathbf{T}_{\varrho,n}\|_2 = \max(4.1 - \varepsilon, 0)$  (magenta dots) in order to stress that the region where Proposition 3.4.4 can be applied is below this barrier.

Figure 3.7, the results of this experiment are shown. From Figure 3.7 (a), we can derive that  $\tilde{q}_{\varepsilon,n}(x) - \min_y \tilde{q}_{\varepsilon,n}(y)$  does nicely peak around the actual support of the ground truth for  $n = 24$  and  $\varepsilon = \|\mathbf{T}_{\varrho,n}\|_2 = 1.5$ . Furthermore, Figure 3.7 (b) indicates that the rate of  $h(\varepsilon^*) \in \mathcal{O}(n^{-3/2})$  seems to hold also for the original objective (3.34) while (c) strongly supports the statement of Proposition 3.4.4 that the choice of  $\varepsilon = \|\mathbf{T}_{\varrho,n}\|_2$  is optimal if

$$\|\mathbf{T}_{\varrho,n}\|_2 = \varepsilon \leq \min\left(\tilde{\sigma}_{|Y|}^{(n)}, \sigma_{|Y|}^{(n)}\right) \leq \min\left(\sigma_3^{(24)} - \|\mathbf{T}_{\varrho,n}\|_2, \sigma_3^{(24)}\right) \approx \min(4.1 - \|\mathbf{T}_{\varrho,n}\|_2, 4.1).$$

In particular, we see that fraction of mass away from the ground truth support increases



for a non-optimal choice of  $\varepsilon$  and we observe that we can naturally only hope for full recovery if  $\|\mathbf{T}_{\varrho,n}\|_2 \leq \sigma_{|Y|}^{(n)}$ .

Even though one might hope to have a simple heuristic parameter choice rule because one can derive  $\varepsilon^* = \|\mathbf{T}_{\varrho,n}\|_2$  without the knowledge of  $\|\mathbf{T}_{\varrho,n}\|_2$  by minimisation of (3.36) with respect to  $\varepsilon$ , this idea is misleading since the computation of (3.36) already includes the knowledge on the underlying support. Therefore, one needs an estimate for the noise level or an experienced practitioner in order to choose  $\varepsilon$  by an a-priori or heuristic parameter choice rule in applications of this method.

---

**Algorithm 1** Support approximation by Christoffel function

---

**Input:** Perturbed moment matrix  $\mathbf{T}_n$ , estimate on contribution  $\|\mathbf{T}_{\varrho,n}\|_2$  of the noise

- 1: Compute low rank approximation  $\mathbf{T}_{\tilde{r},n}$  by truncation of the SVD of  $\mathbf{T}_n$  under the constraint (3.31).
- 2: Choose  $\varepsilon$  according to Proposition 3.4.4.
- 3: Set up the Christoffel function  $\tilde{q}_{\varepsilon,n}$  via (3.32).

**Output:** Visual representation of the support or parameter estimate by computing the local maxima of  $\tilde{q}_{\varepsilon,n}$ .

---

Finally, the method to compute the noisy Christoffel function as a good approximation of the support of a discrete measure is summarised in Algorithm 1. Even if the algorithm does not solve the parameter recovery problem but just gives a visual representation of the measure, the convergence rate from Proposition 3.4.4 is especially remarkable because it gives a completely deterministic result in a natural, multidimensional setting with well-separated nodes and small noise. Moreover, the rate of  $\mathcal{O}(n^{-3/2})$  is faster than the rate of nonlinear polynomial interpolation by the signal polynomial, see Lemma 3.4.2, at the cost of a nonlinear order in the noise level.<sup>74</sup> We postpone more detailed examples for the performance of the method to the next chapter.

---

<sup>74</sup>One should notice that this comparison is slightly flawed because the rates presented in Lemma 3.4.2 and Proposition 3.4.4 are not for the same kind of objective.



## 4 Applications in microscopy

*The second section of this chapter dealing with structured illumination microscopy is related to our publications [70, 69] even though we analyse the problem in a different way in this work where we focus on a condition analysis showing the gain in resolution obtained by this method.*

After observing a limit for the stable recovery of closely spaced positions from low pass Fourier data in Chapter 2, we have seen in Chapter 3 that the presented approximation method performs well if one works above this resolution limit. However, many biological processes proceed on smaller spatial scales such that there has been ongoing research on techniques to overcome the diffraction limit. In this work, we study two of them, namely *Stochastic Optical Reconstruction Microscopy (STORM)* and *Structured Illumination Microscopy (SIM)*, and analyse whether we can understand how these methods overcome the diffraction limit which we defined in Chapter 2.

### 4.1 An approach to STORM analysis

In the late 20th and early 21st century, many approaches to overcome the diffraction limit were proposed in microscopy and culminated into the awarding of the Nobel prize 2014 in Chemistry to Betzig, Hell and Moerner “for the development of super-resolved fluorescence microscopy”, cf. [115]. In particular, highly resolved localisation microscopy where one wants to extract the positions of individual molecules through their light emission is confronted with the problem that on one hand the diffraction limit makes it impossible to recover the positions if the specimen is labelled densely with fluorescent molecules and on the other hand a dense labelling is needed in order to see fine details in the probe. According to [115], the main tool for the solution of this issue was Moerner’s analysis of a labelling molecule which can be turned on and off when it is illuminated by light of a certain wavelength. It was then Betzig who suggested and implemented a first method which used these *photoswitchable* dyes such that out of a dense collection of labels only a small and sparse subset is active at a certain time. This allows to extract the active subset such that one obtains the full information of the dense set of emitters after repeating this process at various time steps. While Betzig published this method as *Photo-Activated Localization Microscopy (PALM)* in 2006<sup>75</sup>, Rust et al. developed *Stochastic Optical Reconstruction Microscopy (STORM)* in the same year and this technique uses basically the same idea with a different dye, see [139]. We visualise this concept of repeated recovery of a stochastically chosen subset of the emitters in Figure 4.1.

Although the strategy behind localisation techniques like STORM appears quite natural, we want to analyse how this method can overcome the diffraction limit as defined in Chapter 2 and we devote the first part of the section to this question. Additionally, it remains a computational challenge to extract the positions of the active molecules in every

---

<sup>75</sup>See [115] for a short overview over the development of super resolution fluorescence microscopy.

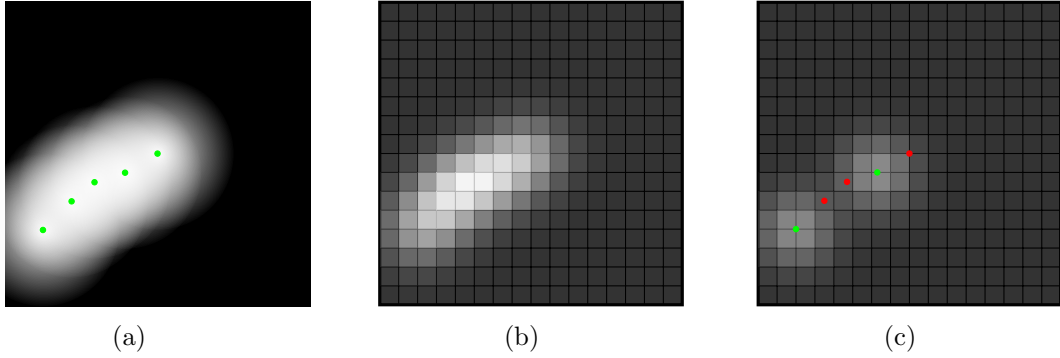


Figure 4.1: The principle of STORM for a toy example. Already for continuous measurements of five blurred, densely spaced labels it is difficult to recover their position (a) and this becomes even worse for realistic discrete data on a grid of pixels in a digital camera (b). On the contrary, turning on only a fraction of the emitters (green) while the rest is off (red) allows to extract the positions of the two active molecules (c). Repeating this stochastic process of photo-switching gives good estimates for the full set of emitter positions.

time step<sup>76</sup> and to bring this information from all frames together in order to obtain a highly resolved image of the specimen. In the second part of this section, we will use the methods from Chapter 3, in particular the signal polynomial, in order to solve this task.

**Gain in condition using STORM** Similar to the previous considerations, we want to recover  $\mu = \sum_{t \in Y} \alpha_t \delta_t$  where  $Y \subset \mathbb{T}^2$  may now contain very closely spaced points compared to the bandlimit  $n$  and the measurements are given by

$$g_s(x) = ([I_s \mu] * h)(x) = \sum_{t \in Y} I_s(t) \alpha_t h(x - t), \quad s = 1, \dots, S \text{ and } x \in [0, 1]^2$$

or equivalently one has access to the Fourier coefficients

$$\widehat{I_s \mu}(k) := \sum_{t \in Y} I_s(t) \alpha_t e^{-2\pi i t k}, \quad s = 1, \dots, S \text{ and } k \in \mathcal{I} = \mathbb{Z}^2 \cap B_n(0). \quad (4.1)$$

Here, the number of frames  $S \in \mathbb{N}$  can typically be of the order  $10^5$  while the most simple way to describe the illumination  $I_s(t)$  is by allowing  $I_s(t) \in \{0, 1\}$  corresponding to the off- and on-state of the photoswitchable dye respectively.<sup>77</sup> As in our analysis of the condition of the original super resolution problem, we mention that previous works can be distinguished between stability analysis for algorithms solving the problem and the study of the condition of the problem itself. A mixture of both approaches can be found in [100] where a situation similar to (4.1) is referred to as “multi snapshot spectral estimation”.<sup>78</sup>

<sup>76</sup>Usually, the image at a fixed time is called frame.

<sup>77</sup>A comparison to [139, Fig. 1b], where the activity of a single label is displayed over time, shows that the assumption  $I_s(t) \in \{0, 1\}$  is reasonable because one can clearly distinguish two different levels of activation alternating over time which correspond either to  $I_s(t) = 0$  or  $I_s(t) = 1$ . Unfortunately, our notation is not optimal at this point because  $t$  runs over the nodes in  $Y$  whereas the natural connotation is that  $t$  represents a variable for time. On the contrary, the variable running through the frames collected over time is denoted by  $s$ . Nevertheless, we keep the notation this way in order to be consistent with the previous chapters.

<sup>78</sup>An important difference is that the weights are modelled as functions of  $s$  and  $t$  without the specific product structure  $I_s(t) \alpha_t$ .

On one hand, the authors derive a variant of MUSIC and ESPRIT for this situation and prove its stability due to the repeated measurements allowing to resolve details finer than the diffraction limit. On the other hand, they show using the Cramer-Rao bound in [100, Thm. VI.3] that each unbiased estimator  $\hat{Y}$  for the node set  $Y$  has a covariance satisfying

$$\mathbb{E} \left[ \text{md}(Y, \hat{Y})^2 \right] \geq c \cdot \frac{(n \cdot \text{sep } Y)^{-2\ell+2} \delta^2}{S n^3 \|\mathbf{X}\|_2} \quad (4.2)$$

where  $\text{md}$  denotes the *matching distance* between two finite sets with equal cardinality,  $c > 0$  is some constant,  $\ell \in \mathbb{N}$  is the size of the largest “clump”, i.e. the largest number of points which are closer spaced than the diffraction limit,  $\delta^2$  is the variance of the assumed Gaussian noise and  $\mathbf{X}$  is a (covariance) matrix containing the weights. It is remarkable that their stability bound for ESPRIT [100, eq. (19)] matches this order in the number of frames  $S$ , the noise level  $\delta$  and the inverse of the *super resolution factor*  $n \cdot \text{sep } Y$ . However, a drawback of this work is its limitation to the univariate case  $Y \subset \mathbb{T}$ .

The work by Liu et al. [107] overcomes this issue and studies the problem (4.1) in the univariate and bivariate case. Their approach combines the various illuminations  $I_s(t)$  as entries of a illumination matrix  $I$  of size  $S \times |Y|$  and quantifies the gain in resolution by examination of the spectral properties of this matrix. More precisely, the minimal separation allowing for approximate recovery of  $\mu$  with the correct number of parameters is

$$\text{sep } Y \geq \frac{c|Y|^2}{n} \left( \frac{\delta}{\|I\|_{\infty, \min} \cdot \alpha_{\min}} \right)^{1/|Y|} \quad (4.3)$$

in the bivariate case if  $c > 0$  is some constant,  $\delta$  the noise level and the influence of the illumination matrix is taken into account by  $\|I\|_{\infty, \min} = \min_{\|x\|_{\infty} \geq 1} \|Ix\|_{\infty} \geq \frac{\sigma_{\min}(I)}{\sqrt{S}}$  (cf. [107, Thm. 3.1]). One should compare this bound with the result from [108, Thm. 2.3 and Prop. 2.4] which we described in (2.3) and this shows how the resolution limit can be decreased by an illumination pattern such that  $\|I\|_{\infty, \min}$  is large. A comparison of (4.3) with (4.2) indicates that the exponent  $1/|Y|$  is not optimal if one includes knowledge about the geometry of the nodes and taking  $1/|Y|$  deteriorates the gain of resolution represented by (4.3) if  $|Y|$  is very large as in STORM.

In order to contribute to these approaches of a theoretical explanation for the resolution of STORM, we develop an argument which should be connected to our results on the resolution limit for classical light microscopy in Chapter 2. The stochastic nature of STORM makes it natural to study the resolution limit by the behaviour of the Cramer-Rao bound for an unbiased estimator of the measure as we did in Subsection 2.2.3. There, we assumed  $S = 1$  and  $I_1(t) = 1$  for all  $t \in Y$  as well as uncorrelated Gaussian noise  $\hat{\rho}$  with variance  $\delta^2$  for the measurements of the moments of  $\mu$ , see (2.27). Here, we will extend this model by allowing each node to be in the on-state ( $I_s(t) = 1$ ) or the off-state ( $I_s(t) = 0$ ) respectively and the choice should be made in a probabilistic manner. The simplest idea is then to introduce a parameter  $p \in [0, 1]$  and to draw  $I_s(t)$  for  $s = 1, \dots, S$ ,  $t \in Y$ , mutually independent from the Bernoulli distribution with parameter  $p$ , i.e.  $\mathbb{P}(I_s(t) = 1) = p$  and the illumination at different frames or different nodes are assumed to be independent.<sup>79</sup> Then, we define the following hierarchical model for STORM.

<sup>79</sup>We remark that especially the assumption of independence is a strong simplification. If we consider [139, Fig. 1b] showing  $I_s(t)$  for fixed  $t$  as a function of  $s$ , this indicates that there is a correlation between consecutive values of  $I_s$ .

#### 4 Applications in microscopy

**Definition 4.1.1** (STORM-model). For each frame, we assume that the random variable  $Z_{s,k}$ ,  $s = 1, \dots, S$ ,  $k \in \mathcal{I} = \mathbb{Z}^2 \cap B_n(0)$ , measuring  $\widehat{I_s \mu}(k)$  follows the probability measure

$$\mathbb{P}(Z \in [a, b] + i[c, d], I_s = \iota_s) = \mathbb{P}\left(\mathcal{CN}(\widehat{I_s \mu}(k), \delta^2 I) \in [a, b] + i[c, d] | I_s = \iota_s\right) \cdot \mathbb{P}(I_s = \iota_s)$$

where for  $\iota_s \in \{0, 1\}^{|Y|}$ ,  $|\iota_s| := \sum_t (\iota_s)_t$ , the last probability is  $\mathbb{P}(I_s = \iota_s) = p^{|\iota_s|} (1-p)^{|Y|-|\iota_s|}$  and the conditional probability follows the complex normal distribution with density

$$f_{s,k}(z) = \frac{1}{2\pi} \exp\left(-\frac{|z - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k}|^2}{2\delta^2}\right)$$

for  $z \in \mathbb{C}$ . Letting the resulting measurements at different frames  $s$  or frequencies  $k$  be independent allows to describe them by a large vector in  $\mathbb{C}^{S \cdot |\mathcal{I}|}$  whose joint density function  $f$  does admit

$$\log f(z, \iota) = \left( -\frac{|z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k}|^2}{2\delta^2} - \log 2\pi + \log\left(p^{|\iota_s|} (1-p)^{|Y|-|\iota_s|}\right) \right)_{k,s} \in \mathbb{C}^{S \cdot |\mathcal{I}|}$$

for  $z \in \mathbb{C}^{S \cdot |\mathcal{I}|}$  and  $\iota_s \in \{0, 1\}^{|Y|}$  for  $s = 1, \dots, S$ .

**Lemma 4.1.2.** *The Fisher information matrix (FIM) of the STORM-model has the form*

$$\begin{aligned} J_{\text{STORM}}(\alpha, Y) &= S \sum_{i=0}^{|Y|-2} \binom{|Y|}{i+2} p^{i+2} (1-p)^{|Y|-2-i} J(\alpha, Y) \\ &+ \frac{\sum_{i=0}^{|Y|-1} \binom{|Y|}{i+1} p^{i+1} (1-p)^{|Y|-1-i} - \sum_{i=0}^{|Y|-2} \binom{|Y|}{i+2} p^{i+2} (1-p)^{|Y|-2-i}}{\delta^2 S^{-1} |\mathcal{I}|^{-1}} \cdot C_\alpha \end{aligned}$$

where  $J(\alpha, Y) \in \mathbb{R}^{3|Y| \times 3|Y|}$  is the FIM for the super resolution model from (2.27) and

$$C_\alpha := \text{diag} \left( \underbrace{1, \dots, 1}_{|Y| \text{ times}}, \underbrace{\frac{\sum_{k \in \mathcal{I}} \|k\|_2^2}{|\mathcal{I}|}}_{\text{twice the absolute value squared weight vector}}, \underbrace{(|\alpha_{t_1}|^2, \dots, |\alpha_{t_{|Y|}}|^2, |\alpha_{t_{|1|}}|^2, \dots, |\alpha_{t_{|Y|}}|^2)}_{\text{twice the absolute value squared weight vector}} \right) \in \mathbb{C}^{3|Y| \times 3|Y|}.$$

*Proof.* The derivatives contained in  $J_{\text{STORM}}(\alpha, Y)$  can be computed as

$$\begin{aligned} \left( \frac{\partial \log f(z, \iota)}{\partial \alpha_{t'}} \right)_{k,s} &= -\delta^{-2} \iota_s(t') \Re \left( e^{2\pi i t' k} \left[ z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right] \right), \\ \left( \frac{\partial \log f(z, \iota)}{\partial (t')_j} \right)_{k,s} &= -\delta^{-2} \iota_s(t') \alpha_{t'} \Re \left( 2\pi i k_j e^{2\pi i t' k} \left[ z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right] \right) \end{aligned}$$

for  $j = 1, 2$ . As in Corollary 2.2.23, the Fisher information matrix has the block structure

$$J_{\text{STORM}}(\alpha, Y) = \mathbb{E}_{Z, I} \left[ \begin{pmatrix} G_1 G_1^* & G_2 G_1^* & G_3 G_1^* \\ G_1 G_2^* & G_2 G_2^* & G_3 G_2^* \\ G_1 G_3^* & G_2 G_3^* & G_3 G_3^* \end{pmatrix} \right]$$

and we investigate the top left block  $G_1 G_1^*$  at first. Entrywise, one obtains

$$\begin{aligned}
 (\mathbb{E}_{Z,I} [G_1 G_1^*])_{t',t''} &= \delta^{-4} \mathbb{E}_{Z,I} \sum_{s,k} \iota_s(t') \iota_s(t'') \Re \left( e^{2\pi i t'' k} \left[ z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right] \right) \\
 &\quad \cdot \Re \left( e^{2\pi i t' k} \left[ z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right] \right) \\
 &= \delta^{-4} \mathbb{E}_{Z,I} \sum_{s,k} \iota_s(t') \iota_s(t'') \left( \frac{1}{2} \cos(2\pi(t'' - t')k) \left| z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right|^2 \right. \\
 &\quad \left. + \frac{1}{2} \Re \left( e^{2\pi i (t'' + t') k} \left[ z_{s,k} - \sum_{t \in Y} \iota_s(t) \alpha_t e^{-2\pi i t k} \right]^2 \right) \right) \\
 &= \delta^{-2} \mathbb{E}_I \sum_{s,k} \iota_s(t') \iota_s(t'') \cos(2\pi(t'' - t')k)
 \end{aligned}$$

because the expectation with respect to  $z$  over the absolute value squared gives  $2\delta^2$  whereas the expectation over the square vanishes.<sup>80</sup> The remaining expectation over  $I$  is by the independence of the illumination at different frames given by

$$\begin{aligned}
 (\mathbb{E}_{Z,I} [G_1 G_1^*])_{t',t''} &= \delta^{-2} \sum_{s,k} \sum_{\iota_s \in \{0,1\}^{|Y|}} p^{|\iota_s|+1} (1-p)^{|Y|-1-|\iota_s|} \iota_s(t') \iota_s(t'') \cos(2\pi(t'' - t')k) \\
 &= \delta^{-2} S \sum_{\iota \in \{0,1\}^{|Y|}} p^{|\iota|+1} (1-p)^{|Y|-1-|\iota|} \iota(t') \iota(t'') \sum_{k \in \mathbb{Z}^2 \cap B_n(0)} e^{-2\pi i (t'' - t') k} \\
 &= \delta^{-2} (\mathcal{A}^* \mathcal{A})_{t',t''} \cdot \begin{cases} \sum_{i=0}^{|Y|-1} \binom{|Y|}{i+1} p^{i+1} (1-p)^{|Y|-1-i} & , t' = t'' , \\ \sum_{i=0}^{|Y|-2} \binom{|Y|}{i+2} p^{i+2} (1-p)^{|Y|-2-i} & , t' \neq t'' , \end{cases}
 \end{aligned}$$

where  $\mathcal{A} \in \mathbb{C}^{|X| \times |Y|}$  is the matrix from Corollary 2.2.23. With this observation at hand, the same calculation can be made for the other blocks of  $J_{\text{STORM}}(\alpha, Y)$  such that we can conclude that the entries of  $J_{\text{STORM}}(\alpha, Y)$  are the entries of  $J(\alpha, Y)$  times a constant depending on whether we deal with an entry on or off the main diagonal of each block. By symmetry, the main diagonals of the off-diagonal blocks are zero.<sup>81</sup> Therefore,  $J_{\text{STORM}}(\alpha, Y)$  and  $J(\alpha, Y)$  are a scalar multiple of each other except for the main diagonal and this leads to the proposed statement.  $\square$

**Lemma 4.1.3.** *We can calculate the sums from the previous Lemma in closed form as*

$$\begin{aligned}
 \sum_{i=0}^{|Y|-2} \binom{|Y|}{i+2} p^{i+2} (1-p)^{|Y|-2-i} &= 1 - (1-p)^{|Y|} - |Y| p (1-p)^{|Y|-1}, \\
 \sum_{i=0}^{|Y|-1} \binom{|Y|}{i+1} p^{i+1} (1-p)^{|Y|-1-i} &= 1 - (1-p)^{|Y|}.
 \end{aligned}$$

*Proof.* The binomial coefficients satisfy

$$\binom{|Y|}{i+2} = \binom{|Y|-2}{i} \frac{(|Y|-1)|Y|}{(i+2)(i+1)} \quad \text{and} \quad \binom{|Y|}{i+1} = \binom{|Y|-1}{i} \frac{|Y|}{i+1}$$

<sup>80</sup>This is by definition of the complex normal distribution having relation matrix or pseudo-covariance matrix equal to zero.

<sup>81</sup>The main diagonals have the form  $\alpha_t^2 \sum_{k \in \mathbb{Z}^2 \cap B_n(0)} (2\pi i k_1)^{j_1} (2\pi i k_2)^{j_2}$  for  $(j_1, j_2) \in \{(0,1), (1,0), (1,1)\}$ .

#### 4 Applications in microscopy

such that the result follows by the representation through the integrals  $\frac{p^{i+1}}{i+1} = \int_0^p x^i dx$  and  $\frac{p^{i+2}}{(i+1)(i+2)} = \int_0^p \int_0^x y^i dy dx$  together with the binomial theorem.  $\square$

**Theorem 4.1.4** (Condition of STORM). *The STORM model admits*

$$J_{\text{STORM}}(\alpha, Y) = S \left( 1 - (1-p)^{|Y|} - |Y|p(1-p)^{|Y|-1} \right) J(\alpha, Y) \\ + S|Z|\delta^{-2}|Y|p(1-p)^{|Y|-1}C_\alpha$$

and thus the STORM model is well-conditioned in the sense of Definition 2.2.24 for illumination parameter  $p \in (0, 1)$  regardless of the separation of  $Y$ .

*Proof.* The result for  $J_{\text{STORM}}(\alpha, Y)$  follows by Lemma 4.1.2 and Lemma 4.1.3 such that the smallest eigenvalue of it can be estimated as

$$|Z|^{-1}\lambda_{\min}(J_{\text{STORM}}(\alpha, Y)) \geq \delta^{-2}S|Y|p(1-p)^{|Y|-1} \min(1, \alpha_{\min}^2). \quad (4.4)$$

As this lower bound is independent of  $n$  and positive for  $p \in (0, 1)$ , Definition 2.2.24 is fulfilled for any separation of the node set  $Y$ .  $\square$

Apart from the natural observation that  $J_{\text{STORM}}(\alpha, Y)$  equals  $J(\alpha, Y)$  for  $S = 1$  and  $p = 1$ , we complete the analysis of the resolution of STORM with the following remarks.

**Remark 4.1.5.** (i) First of all, it makes sense that the lower bound (4.4) deteriorates for  $p \in \{0, 1\}$  because one cannot hope for an improved condition if no labels or all labels emit light at the same time. An intuition for the parameter  $p$  might be given by the observation that  $|Y| \cdot p$  describes the expected number of active labels in each frame. Setting this value to  $r$ , the lower bound on  $|Z|^{-1}\lambda_{\min}(J_{\text{STORM}}(\alpha, Y))$  from (4.4) is

$$\delta^{-2}Sr \left( 1 - \frac{r}{|Y|} \right)^{|Y|-1} \min(1, \alpha_{\min}^2) \approx \delta^{-2}Sr \left( 1 - \frac{r}{|Y|} \right)^{-1} e^{-r} \min(1, \alpha_{\min}^2)$$

if we think of  $|Y|$  being large. Hence, we see that the expected number of on-labels per frame needs to be nicely balanced in order to obtain a large value for the lower bound in (4.4).

- (ii) As for Vandermonde matrices, there is some theory available for the smallest singular value of a confluent Vandermonde matrix with clustering, univariate nodes, see [9, 100], showing that  $\lambda_{\min}(J(\alpha, Y))$  behaves like  $(n \cdot \text{sep } Y)^{4\ell-2}$  if  $\ell \in \mathbb{N}$  is the size of the largest cluster, e.g. cf. [9, Thm. 3.1], and we might conjecture a similar behaviour for higher dimensions. Under the assumption of such a conjecture, one could then obtain a bound similar to (4.2) with the same dependency on  $S$  but with a different behaviour in the super resolution factor  $(n \cdot \text{sep } Y)^{-1}$ . The latter might be because (4.2) takes into account estimates for  $Y$  only whereas we consider estimators for  $Y$  and the weights simultaneously.
- (iii) Although we already mentioned limitations of the model, it is justified as it is able to describe the gain of resolution through STORM and we want to highlight that it especially reflects the intuition behind the parameter  $S$ . For example, it is natural that the model gives  $J_{\text{STORM}}(\alpha, Y) = S \cdot J(\alpha, Y)$  for  $p = 1$  because the  $S$ -fold, independent repetition of the original super resolution problem would lead to the same amount of reduction in the variance of the noise, see [100, p. 4564].



**Application of Christoffel function and signal polynomial to STORM data** A wide range of algorithms for the processing of STORM data is available. Most common might be algorithms which use filters to detect the parts of each frame where molecules are existing and then find their sub-pixel positions by fitting a PSF into each of the regions. For example, *ThunderSTORM* [125] belongs to this class of methods. In contrast to this, mathematically more sophisticated algorithms using SDP optimisation or gradient methods like the sliding Frank-Wolfe method (e.g. [17, 96, 35]) or (multi-snapshot) subspace methods (e.g. Li et al. [100]) appear to be used less frequently. Instead, recent approaches involving the utilisation of deep neural networks have gained some attention, see [122, 121].

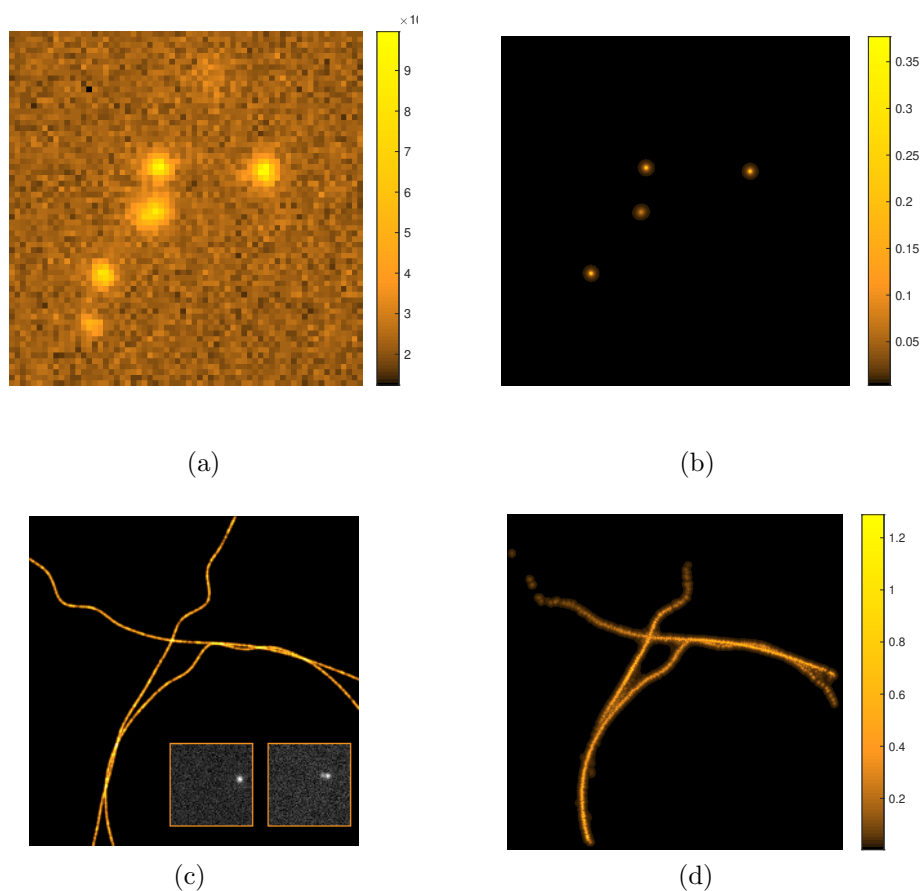


Figure 4.2: STORM data set and results. An individual frame (a) and its reconstruction via  $(1 - \tilde{p}_{1,n})^{-1}$  (b) is shown in the first row. In the second row, the ground truth containing a simulation of *microtubules* (c) taken from <https://srm.epfl.ch/Challenge> is compared with our result obtained by taking the pixelwise maximum of  $(1 - \tilde{p}_{1,n})^{-1}$  over all frames (d).

In order to test our approach from Chapter 3 to interpolate the support of a measure by the signal polynomial or the Christoffel function, we analyse the test data set MTO.N1.HD which is publicly available in the context of the EPFL SMLM challenge 2016, see <https://srm.epfl.ch/Challenge>. It consists of a stack of 2500 STORM images of size  $64 \times 64$  simulated from a ground truth distribution of fluorophores depicted in Figure 4.2 (c). Exemplarily, the tenth frame is shown in Figure 4.2 (a) in order to give an idea on the *signal-to-noise-ratio* (SNR) and the number of active molecules in each frame.

As for the previous numerical experiments, the code can be found on <https://github.com/MHockmann/Dissertation.git>. At first and for each frame, one needs to compute the moments of the underlying discrete measure such that deconvolving with a carefully estimated PSF and thresholding the considered range of moments at some frequency  $n$  are needed. The theoretical resolution limit gives a first idea how to choose the frequency parameter  $n$ . While the size of the field of view in  $x$  and  $y$  direction is simulated to be 6400 nm and the wavelength of the emitted light is set to 660 nm, the simulation uses a *numerical aperture* (NA) of 1.49 such that the theoretical bandlimit is expected<sup>82</sup> to be at  $1.49 \cdot 2 \cdot \frac{6400}{660} \approx 29$ . However, one has to make sure that the moments can be distinguished from the noise such that we consider the moments  $\hat{\mu}(k)$  with frequencies  $k \in \mathbb{Z}^2 \cap B_n(0)$  and  $n = 20$ . For these frequencies, the moments are obtained with a Gaussian PSF whose standard deviation is heuristically chosen such that the first image is nicely deconvolved.

The computation of the signal polynomial and of the Christoffel function demand for the SVD of the moment matrix. As it is already costly to store this matrix of size  $|\mathcal{I}| \times |\mathcal{I}|$  with  $|\mathcal{I}| = 317 \approx \pi(n/2)^2$  for each frame, we just implement the matrix-vector multiplication for  $\mathbf{T}_n$  with any vector  $v$  by observing that

$$(\mathbf{T}_n v)_\ell = \int_{\mathbb{T}^2} e^{-2\pi i \ell x} \left( \sum_{k \in \mathbb{Z}^2 \cap B_{n/2}(0)} e^{2\pi i k x} v_k \right) d\mu(x), \quad \ell \in \mathbb{Z}^2 \cap B_{n/2}(0), \quad (4.5)$$

can be approximated efficiently by two FFTs without storage of  $\mathbf{T}_n$ . This allows to calculate only the first singular vectors and values in a sparse SVD by the MATLAB function `svds`.

---

**Algorithm 2** STORM processing with signal polynomial

---

**Input:** Images  $g_s$ ,  $s = 1, \dots, S$ , PSF  $h$  and bandlimit  $n$

- 1: Obtain low order moments up to order  $n$  by deconvolution, set  $s = 1$
  - 2: **while**  $s \leq S$  **do**
  - 3:   Compute sparse SVD of  $\mathbf{T}_n$  by implementing (4.5) with FFTs
  - 4:   Truncate SVD at singular values below certain fraction of the largest singular value
  - 5:   Store evaluation of signal polynomial  $\tilde{p}_{1,n,s}$
  - 6:    $s = s + 1$
  - 7: **end while**
  - 8: For each  $x$  compute  $\max_s (1 - \tilde{p}_{1,n,s}(x))^{-1}$
- 

It turns out that the influence of the noise is larger than assumed in Section 3.4 such that only the singular vectors corresponding to the largest singular values are not dominated by noise. Hence, the regularisation parameter  $\varepsilon$  in  $\tilde{q}_\varepsilon$  would need to be too large in order to guarantee a well-localised interpolation of the support with the Christoffel function. Additionally, the reconstructed intensities are not only secondary but also harm the reconstruction of the full image because merging images with possibly very different intensities is problematic. Consequently, we refuse to consider the Christoffel function but use the interpolation property of signal polynomial, see Section 3.2. Its noisy version  $\tilde{p}_{1,n}(x)$  is large (with value slightly smaller than 1) if  $x$  is close to the support such that we plot  $(1 - \tilde{p}_{1,n})^{-1}$  for the representation of the measure.<sup>83</sup> In fact, this is then similar

<sup>82</sup>By Abbe's diffraction limit, the inverse of the bandlimit  $n$  admits  $n^{-1} = \frac{\lambda}{2NA}$  if  $\lambda$  is the wavelength of the emission light and taking into account the size of the field of view gives the mentioned estimate.

<sup>83</sup>We evaluate this super-resolved representation on a grid of  $2048 \times 2048$  pixels.

to MUSIC, cf. [144]. It can be seen in Figure 4.2 that this yields very highly resolved representations of the expected ground truth even though choosing the numerical rank for  $\tilde{p}_{1,n}$  too small might lead to false negative identifications, i.e. there is the risk to identify only four out of five active molecules in this frame. Nevertheless, balancing the risk of a few false negatives in a stack of 2500 frames with the risk of false positive identifications gives rise to a small value for the number of terms included in the signal polynomial.<sup>84</sup>

The ratio that each support identification should be represented in the overall approximation motivates to take for each pixel the maximum of  $(1 - \tilde{p}_{1,n})^{-1}$  over all frames and the result is shown in Figure 4.2 (d).<sup>85</sup> The outcome has a quality similar to the reconstruction in [96] showing in particular at least in outlines the interweaving microtubules on the right part of the image. Our approach which we summarise in Algorithm 2 needs less than a second per frame such that parallelisation enables to solve the task within minutes and by improving the implementation further speed-up might be possible. However, it appears not to be realistic to exceed the performance of an approach based on deep learning like DeepSTORM, see [122, 121], with this kind of method.

## 4.2 Structured Illumination Microscopy

The idea to enhance the resolution of an optical system by manipulation of the illumination of the sample is used not only for stochastic illumination as in STORM but also for illumination with deterministic patterns. An example for the latter is *Structured Illumination Microscopy (SIM)* introduced by Heintzmann [65, 66] and Gustafsson [62] where typically a periodic, nonnegative illumination function

$$I(x) = \sum_{m=-M}^M b_m e^{2\pi i m(vx + \varphi)}, \quad b_m \in \mathbb{C}, \overline{b_m} = b_{-m}, v \in \mathbb{R}^2, \varphi \in \mathbb{T}, M \in \mathbb{N},$$

is used in order to obtain measurements

$$g(x) = ([I\mu] * h)(x) = \sum_{t \in Y} I(t) \alpha_t h(x - t)$$

at discrete values  $x = x_j = \frac{j}{J}, j \in \{0, \dots, J-1\}^2$  and sampling parameter  $J \in \mathbb{N}$ .<sup>86</sup> Here, the situation where  $M = 1$  and thus  $I(x) = 1 + c_0 \cos(2\pi(vx + \varphi))$  is called *linear SIM* with *modulation depth*  $c_0 \in (0, 1)$ , while an illumination function with order  $M > 1$  can be generated for instance by utilisation of saturation effects and the resulting microscopy technique is then called *nonlinear SIM*.<sup>87</sup> The main idea of SIM is that modulation in real space corresponds to translations in Fourier space such that by multiplication of the ground truth  $\mu$  with  $I$  the spectral data contains a collection of translates, i.e.

$$\hat{g}(k) = \sum_{m=-M}^M b_m e^{2\pi i m \varphi} \hat{h}(k) \cdot \hat{\mu}(k - mv), \quad k \in \mathbb{Z}^2.$$

<sup>84</sup>We set the maximal number of considered singular vectors to eight.

<sup>85</sup>A more sophisticated statistical approach would be to take the second or third highest value as this might reduce the probability of a false positive support identification. Taking the maximum value is often called *maximum intensity projection* in volumetric imaging, e.g. see [153].

<sup>86</sup>Compare this to our setting in Subsection 2.2.4.

<sup>87</sup>For example, nonlinear SIM was introduced in [66, 63]. See [75] for a good overview over linear and nonlinear SIM.

Despite the fact that  $\hat{h}$  is typically bandlimited, the signal  $g$  thus contains frequency information of  $\hat{\mu}$  outside of  $\text{supp } \hat{h}$ . Given the knowledge of the optical transfer function (OTF)  $\hat{h}$  and measurements  $g_{l,s}$  for various angles  $\varphi_s$  and pattern vectors  $v_l$  it is then possible to recover the spectrum of  $\mu$  on a larger region in Fourier space. Schematically, this is displayed in Figure 4.3 for the classical case of three shift vectors  $v_l, l = 1, 2, 3$ . The method to process the data is referred to as the *Gustafsson algorithm* and it is well described in [75, 92]. However, we observed in [69] that this Gustafsson algorithm may

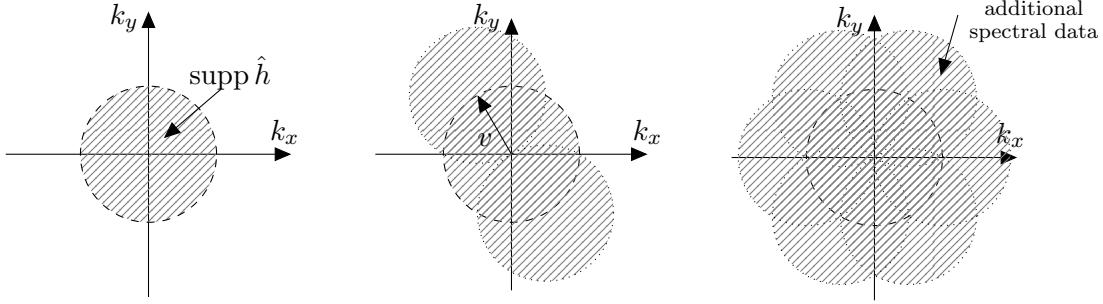


Figure 4.3: Linear 2D-SIM in frequency space: While conventional microscopy just allows to reconstruct the spectrum  $\hat{\mu}$  on  $\text{supp } \hat{h}$  (left), SIM-images contain shifted spectral data in directions  $\pm v$  (middle). By rotation of the pattern one obtains additional data in various directions (right).

fail if the ground truth  $\mu$  is a discrete measure.<sup>88</sup> Additionally, the expression  $\hat{\mu}(k - mv)$  for  $k \in \mathbb{Z}^2, m \neq 0, v \notin \mathbb{Z}^2$ , makes only sense in the non-periodic setting. Hence, we analyse the problem in the setting from Subsection 2.2.4 and study how the gain in resolution through SIM can be justified by our model of a diffraction limit.<sup>89</sup> Furthermore, we want to find a way to circumvent the issue with the Gustafsson algorithm for discrete measures.

**Resolution analysis for SIM** Let the PSF  $h$  admit the properties of Definition 2.2.29. We generalise the model from Subsection 2.2.4 and assume access to discrete samples of

$$\begin{aligned} g_{l,s}(x) &= [h * (I_{s,l}\mu)](x) = \left[ h * \left( \sum_{m=-1}^1 b_m e^{2\pi i \cdot m v_l + 2\pi i m \varphi_s} \mu \right) \right](x) \\ &= \sum_{m=-1}^1 e^{2\pi i m \varphi_s} [h * (b_m e^{2\pi i \cdot m v_l} \mu)](x) \end{aligned} \quad (4.6)$$

at all points  $x = x_j = \frac{j}{J} := (j_1/J, j_2/J)^\top \in [-\frac{1}{2} + \Delta, \frac{1}{2} + \Delta]^2$  for  $j_1, j_2 \in \mathbb{Z}$  with field-of-view parameter  $\Delta > 0$ , angles  $\varphi_s, s = 1, \dots, 3$ , and three modulation vectors  $v_l$  with indices  $l = 1, 2, 3$ . We remark that this is the classical linear SIM setting of three orientations and that our approach can be extended to nonlinear SIM or to different illumination patterns,

<sup>88</sup>Summarised briefly, the problem is that the decoupled data for  $\hat{\mu}(k + mv)$  is given on an integer grid  $k \in \mathbb{Z}^d$  while the translation vector  $v$  generically satisfies  $v \notin \mathbb{Z}^2$ . This means that an interpolation is necessary in order to evaluate all decoupled data sets on the same grid. Therefore, the data is shifted in Fourier space by multiplication with an appropriate modulation in real space but this works well only for data with few oscillations in Fourier domain which is not the case for discrete measures.

<sup>89</sup>As  $\|v\|_2 \leq n$  where  $n$  is the bandlimit of  $h$ , see [75, p. 5], linear SIM is often advertised to increase the resolution up to a factor two while nonlinear SIM taking more translates into account can gain even better resolution (e.g. cf. [62, 63]).

e.g. to the ones presented by Ingerman et al. in [75]. Moreover, we expect the angles  $\varphi_s$  and pattern vectors  $v_l$  to be known.<sup>90</sup> Then, the classical step of the Gustafsson method would be to compute (an approximation to) the Fourier transform of the data and to disentangle the contribution of the different orders in  $m$  afterwards, cf. [75, 92]. However, the signal is much more regular (in the sense that it contains only a few oscillations) in spatial domain than in Fourier domain. Therefore, we propose to separate the components before taking the Fourier transform. This will then also allow to circumvent the issues with the Gustafsson algorithm in the next paragraph. For the analysis of the condition, we compute

$$f_{l,m}(x_j) := \frac{1}{3} \sum_{s=1}^3 e^{-2\pi i m \varphi_s} g_{l,s}(x_j) \quad (4.7)$$

and obtain the bandlimited function values

$$\begin{aligned} f_{l,m}(x_j) &= [h * (b_m e^{2\pi i m v_l} \mu)](x_j) = b_m \sum_t \alpha_t e^{2\pi i m v_l t} h(x_j - t) \\ &= b_m \sum_t \alpha_t \int_{\mathbb{R}^2} \hat{h}(\xi) e^{-2\pi i t(\xi - m v_l)} e^{2\pi i \xi x_j} d\xi \\ &= e^{2\pi i m v_l x_j} \int_{\mathbb{R}^2} b_m \hat{h}(\xi + m v_l) \hat{\mu}(\xi) e^{2\pi i \xi x_j} d\xi \end{aligned} \quad (4.8)$$

for  $m = -1, 0, 1$  if the angles  $\varphi_s, s = 1, 2, 3$ , are equidistant in  $\mathbb{T}$  such that (4.7) is the application of the inverse of the Fourier matrix  $(e^{2\pi i m \varphi_s})_{1 \leq s, m \leq 3}$ .<sup>91</sup> As in Subsection 2.2.4, we denote the vectors of evaluations of  $g_{l,s}$  and  $f_{l,m}$  at  $x_j$  by  $\tilde{g}_{l,s}$  or  $\tilde{f}_{l,m}$  respectively. This allows to generalise the definition of the reconstruction map from Definition 2.2.30.

**Definition 4.2.1** (SIM-Reconstruction from image data). For given PSF  $h$  as in Definition 2.2.29, illuminations  $I_{s,l}, 1 \leq l, s \leq 3$  as in (4.6), and  $q, \Delta > 0$  the linear *SIM-image data reconstruction map* is  $\tilde{\mathcal{R}}_{\text{SIM}} : \mathbb{C}^{9\mathcal{J}} \rightarrow \mathcal{P}(\mathcal{M}(q))$ ,

$$\tilde{g} \mapsto \operatorname{argmin}_{\nu \in \mathcal{M}(q)} \|\tilde{g} - h * [I\nu]\|_{\text{SIM},2}^2 := \operatorname{argmin}_{\nu \in \mathcal{M}(q)} \sum_{s,l} \sum_{x_j \in [-\frac{1}{2} + \Delta, \frac{1}{2} + \Delta]^2} \left| (\tilde{g}_{s,l})_j - (h * [I_{s,l}\nu])(x_j) \right|^2$$

where we consider in this section  $\mathcal{M}(q)$  as the set of *non-periodic*, discrete measures with support  $Y \subset [-\frac{1}{2}, \frac{1}{2}]^2$  having *Euclidean* separation  $\min_{t, t' \in Y} \|t - t'\|_2$  at least  $q$ .

**Definition 4.2.2** (Condition for SIM). We define the *condition number of SIM* as<sup>92</sup>

$$\tilde{\kappa}_{\text{SIM}}(q, \Delta, J, h, M) := \sup_{\substack{\mu \in \mathcal{M}(q) \\ |Y^\mu| \leq M}} \sup_{\substack{\rho \in \mathbb{C}^{9\mathcal{J}} \\ \rho \neq 0}} \inf_{\nu \in \mathcal{R}_{\text{SIM}}((h * [I\mu])(x_j))_j + \rho)} \frac{W_1(\nu, \mu)}{\|\rho\|_{\text{SIM},2}}.$$

The condition can then be analysed by the Cauchy-Schwartz inequality in the estimation

$$\sum_{m,l} \sum_{x_j \in [-\frac{1}{2} + \Delta, \frac{1}{2} + \Delta]^2} \frac{|g_{l,m}(x_j)|^2}{3} = \sum_{m,l} \sum_{x_j \in [-\frac{1}{2} + \Delta, \frac{1}{2} + \Delta]^2} |f_{l,m}(x_j)|^2 \geq \sum_{x_j \in [-\frac{1}{2} + \Delta, \frac{1}{2} + \Delta]^2} \frac{|\sum_{m,l} |f_{l,m}|^2}{9}$$

<sup>90</sup>How these parameters can be derived from the data is explained in [92].

<sup>91</sup>Note that this matrix is well conditioned such that the separation of the components does not affect the analysis of the condition. Of course, the latter is also the motivation to achieve equidistant angles in practice.

<sup>92</sup>As in Subsection 2.2.4, we use the 1-Wasserstein distance according to Proposition 1.4.5 with  $\mathcal{X} = [-\frac{1}{2}, \frac{1}{2}]^2$  equipped with the Euclidean distance. The norm  $\|\cdot\|_{\text{SIM},2}$  was introduced in Definition 4.2.1.

#### 4 Applications in microscopy

and thus

$$\begin{aligned} \sum_{m,l} \sum_{x_j \in [-\frac{1}{2}+\Delta, \frac{1}{2}+\Delta]^2} |g_{l,m}(x_j)|^2 &\geq \sum_{x_j \in [-\frac{1}{2}+\Delta, \frac{1}{2}+\Delta]^2} \left| \int_{\mathbb{R}^2} \left( \sum_{m,l} b_m \hat{h}(\xi + mv_l) \right) \frac{\hat{\mu}(\xi)}{\sqrt{3}} e^{2\pi i \xi x_j} d\xi \right|^2 \\ &= \frac{1}{3} \sum_{x_j \in [-\frac{1}{2}+\Delta, \frac{1}{2}+\Delta]^2} |(h_{\text{SIM}} * \mu)(x_j)|^2 \end{aligned}$$

with the *SIM-PSF*  $h_{\text{SIM}}(x) = \sum_{m,l} b_m e^{2\pi i m v_l x} h(x)$  satisfying  $\text{supp } \hat{h}_{\text{SIM}} = \bigcup_{m,l} B_n(mv_l)$ . With this extension of the support of the OTF, it is straightforward to see that the following definition of the diffraction limit for SIM being analogously to Definition 2.2.35 can be made precise as a corollary of Theorem 2.2.36.

**Definition 4.2.3** (Diffraction limit of SIM). We define the *optimal transition constant of SIM* as

$$\tilde{\Omega}_{\text{SIM},2} = \inf \left\{ \tilde{q} > 0 : \exists \beta \in \mathbb{N} \lim_{n \rightarrow \infty} \sup_{M \leq (\sqrt{dn}/\tilde{q})^d} \frac{\tilde{\kappa}_{\text{SIM}} \left( \frac{\tilde{q}}{n}, \Delta, J, h, M \right)}{M^\beta} < \infty \right\}.$$

**Corollary 4.2.4.** Let the wave vectors  $v_l$  satisfy  $\|v_l\|_2 = k_0 \in [0, n]$ . Then, we have

$$\frac{1}{n+k_0} \sqrt{\frac{4}{3}} \leq \tilde{\Omega}_{\text{SIM},2} \leq \frac{1}{\frac{\sqrt{3}}{2}k_0 + \sqrt{n^2 - \frac{1}{4}k_0^2}} \frac{j_{1,1}}{\pi}$$

for the optimal transition constant of linear SIM.

*Proof.* Because of  $h_{\text{SIM}}(x) = I(x) \cdot h(x)$  and  $k_0 \leq n$ , the estimates on the polynomial decay of  $h$  and its radial derivative carry over to the ones for  $h_{\text{SIM}}$ . In addition, the support of  $\hat{h}_{\text{SIM}}$  admits

$$B_{\frac{\sqrt{3}}{2}k_0 + \sqrt{n^2 - \frac{1}{4}k_0^2}}(0) \subset \text{supp } \hat{h}_{\text{SIM}} = \bigcup_{m,l} B_n(mv_l) \subset B_{n+k_0}(0)$$

and this allows to conclude the statement by Theorem 2.2.36.  $\square$

This analysis shows nicely that our definition of the diffraction limit is also able to explain the increased resolution by SIM. Finally, we remark that this result could similarly be extended to nonlinear SIM. In this case, one would need to know the behaviour of the coefficients  $b_m$  and this is beyond the scope of this work.

**SIM algorithm for discrete measures** The approach to separate the spectral components in the spatial domain as described in (4.8) is not only helpful for the analysis of the condition of SIM but also beneficial for the derivation of a SIM-algorithm for sparse measures. From (4.8), we can compute

$$\begin{aligned} \sum_{m,l} e^{-2\pi i m v_l x_j} f_{m,l}(x_j) &= \int_{\mathbb{R}^2} \left( \sum_{m,l} b_m \hat{h}(\xi + mv_l) \right) \hat{\mu}(\xi) e^{2\pi i \xi x_j} d\xi \\ &= \left[ \left( \sum_{m,l} b_m e^{-2\pi i m v_l y} h(y) \right) * (\mu(y)) \right] (x_j) = (h_{\text{SIM}} * \mu)(x_j) \end{aligned}$$

and thus one obtains a bandlimited version of  $\mu$  with larger bandlimit than the original data  $g_{l,s}$ . From there it is then straightforward to obtain  $\hat{\mu}$  on a larger grid of frequencies, which allows the estimation of finer details. We summarise this in Algorithm 3. As an alternative to the last step, one might use Chapter 3 in order to obtain an approximate representation of  $\mu$  by a polynomial or a rational function.

---

**Algorithm 3** SIM with separation and recombination in spatial domain.

---

**Input:**  $g_{l,s}(x_j)$

- 1: Separate components  $f_{l,m}$  by DFT along the arc's  $\varphi_s$  as in (4.7).
  - 2: Compute  $\sum_{m,l} e^{-2\pi i m v_l x_j} f_{m,l}(x_j)$ .
  - 3: Approximate its Fourier transform by FFT.
  - 4: Obtain approximate values for  $\hat{\mu}$  on extended, equispaced grid by deconvolution.
  - 5: Use a standard algorithm like Matrix Pencil, MUSIC or ESPRIT in order to estimate parameters of  $\mu$ .
- 

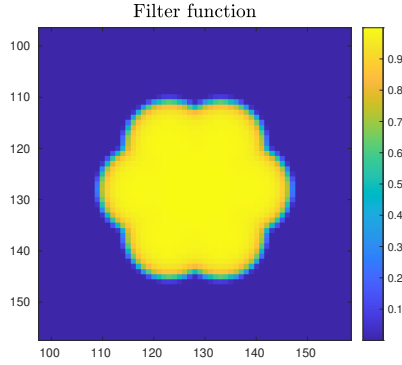


Figure 4.4: Low-pass filter given by the fraction in (4.10) for the experiment in Example 4.2.6 where  $n = k_0 = 10$  and  $\theta = 10^{-3}$ .

In contrast to Algorithm 3 one could also perform the recombination in Fourier domain after doing the order separation in spatial domain. One would then compute an estimate for the Fourier transform of  $e^{-2\pi i m v_l} f(\cdot)$  followed by a *Wiener filter* in Fourier space.<sup>93</sup> More precisely, one obtains the Wiener estimate

$$\hat{\mu}_\theta(\xi) := \frac{\sum_{m,l} b_m \hat{h}(\xi + m v_l) [e^{-2\pi i m v_l} f_{m,l}(\cdot)](\xi)}{\sum_{m,l} |b_m|^2 |\hat{h}(\xi + m v_l)|^2 + \theta} \quad (4.9)$$

for some regularisation parameter  $\theta > 0$ . By substituting (4.8), this reads

$$\hat{\mu}_\theta(\xi) = \frac{\sum_{m,l} |b_m|^2 |\hat{h}(\xi + m v_l)|^2}{\sum_{m,l} |b_m|^2 |\hat{h}(\xi + m v_l)|^2 + \theta} \cdot \hat{\mu}(\xi) \quad (4.10)$$

in the noise free case such that  $\theta = 0$  would lead to the perfect low-pass filter in this situation. For a small but positive value of  $\theta$ , the filter is displayed in Figure 4.4. To sum up, the approach with the Wiener filter for recombination in Fourier domain is then outlined in Algorithm 4.

---

<sup>93</sup>The Wiener filtering in Fourier domain is also part of the classical Gustafsson method, cf.[75, 92]. Its background is that one is confronted with different estimates for the same moments of  $\mu$  at frequencies where the supports of the translated OTF overlap, see Figure 4.3. Wiener filtering allows to combine these measurements by weighting each estimate by the value of the translated OTF at this frequency. As described in [75], this is the optimal combination with respect to the noise.

**Algorithm 4** SIM with separation in spatial and recombination in Fourier domain**Input:**  $g_{l,s}(x_j)$ , regularisation parameter  $\theta$ 

- 1: Separate components  $f_{l,m}$  by DFT along the arc's  $\varphi_s$  as in (4.7).
- 2: Compute the Fourier transform  $[e^{-2\pi i m k_l} \cdot f_{m,l}(\cdot)]^\wedge$ .
- 3: Obtain a noise-optimal estimate for  $\hat{\mu}_\theta$  on enlarged set by Wiener filter (4.9).
- 4: Estimate parameters of  $\mu$  or approximate  $\mu$  by polynomial or rational function.

**Remark 4.2.5.** Algorithm 4 has the advantage that the recombination (4.9) is provably noise-optimal in Fourier domain whereas it comes with the drawbacks that we have to compute more FFTs compared to Algorithm 3 and that we have to choose the regularisation parameter  $\theta$  reasonably. One might define parameter choice rules by studying the connection of (4.9) to a Tikhonov-regularised least squares problem. For those problems, a large range of parameter choice rules is available. On the other hand, one could also modify Algorithm 3 by allowing additional weights  $\beta_{m,l}$  in the recombination  $\sum_{m,l} \beta_{m,l} e^{-2\pi i m v_l x_j} f_{m,l}(x_j)$  and optimising them subject to a desired criterion. However, this is beyond the scope of this work.

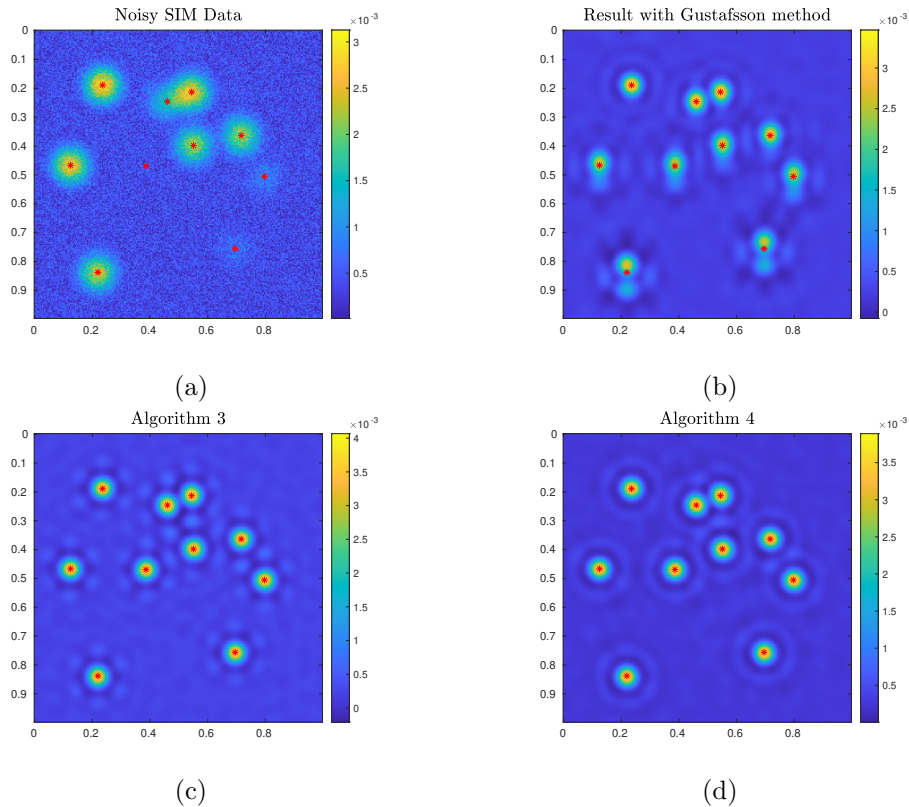


Figure 4.5: Results of Algorithm 3 and Algorithm 4 compared to Gustafsson algorithm for simulated data consisting of ten nodes having uniform weights. For  $n = k_0 = 10$ , we simulate noisy SIM images  $g_{l,s}$  for three pattern directions  $v_l$  and three angles  $\varphi_s$  for each direction (a). From this, the result of the Gustafsson method containing artefacts (b) as well as the results of the presented algorithms (c,d) are shown. In (b),(c),(d) we represent the moments as the outcome of the algorithms by their filtering through an *apodisation function*.



**Example 4.2.6.** In Figure 4.5, we compare the performances of Algorithm 3 and Algorithm 4 with the outcome of the Gustafsson method in a simple example consisting of a measure with ten randomly selected nodes and weights uniformly equal to one. The simulated SIM images for  $n = k_0 = 10$  and coefficient vector  $b = (\frac{1}{2}, 1, \frac{1}{2})^\top$  contain noise where the noise is drawn from a uniform distribution with maximal intensity equal to 30% of the largest intensity of the signal. In order to compare the moments resulting as estimates for  $\hat{\mu}$  from the algorithms and as a way to reduce the ringing, we filter these values with a window or *apodisation function*  $a$  and plot

$$\sum_{k \in \mathbb{Z}^2} \hat{a}(k) \hat{\mu}(k) e^{2\pi i k x}$$

via the FFT for each of the estimates  $\hat{\mu}$ . Therein, the apodisation function is chosen such that  $\hat{a}$  goes to zero smoothly next to the maximal observed frequency and more precisely

$$\hat{a}(k) = \exp\left(\left(1.9n\right)^{-1} - \left(0.1(1.9 \cdot n - \|k\|_2^2)\right)^{-1}\right)$$

for  $\|k\|_2 \leq 1.9 \cdot n$  and zero else.<sup>94</sup> As described before, the Gustafsson methods fails to give a meaningful result because we use pattern vectors  $v_l \notin \mathbb{Z}^2$ . In contrast to this, Algorithm 3 and Algorithm 4 provide higher resolved representations of the measure compared to the SIM data such that a gain in resolution is indeed observable. In this example, Algorithm 4 seems to produce less pattern artefacts than Algorithm 3 even though the differences appear to be marginal.

---

<sup>94</sup>Simpler choices like triangular windows with less regularity are possible as well, e.g. cf. [75].



## Bibliography

- [1] E. Abbe. Beiträge zur Theorie des Mikroskops und der mikroskopischen Wahrnehmung. *Arch. Mikr. Anat.*, 9:413–468, 1873.
- [2] G. B. Airy. On the diffraction of an object-glass with circular aperture. *Trans. Cambridge Philos.*, 5:283–291, 1835.
- [3] N. Akhiezer and M. Krein. On the best approximation of periodic functions. *Dokl Akad. Nauk SSSR*, 15:107–112, 1937.
- [4] F. Andersson and M. Carlsson. ESPRIT for multidimensional general grids. *SIAM J. Mat. Anal. Appl.*, 39(3):1470–1488, 2018.
- [5] K. Atkinson and W. Han. *Theoretical numerical analysis*, volume 39 of *Texts in Applied Mathematics*. Springer, Dordrecht, third edition, 2009.
- [6] C. Aubel and H. Bölcskei. Vandermonde matrices with nodes in the unit disk and the large sieve. *Appl. Comput. Harmon. Anal.*, 47(1):53–86, 2019.
- [7] C. Aubel and H. Bölcskei. Deterministic performance analysis of subspace methods for cisoid parameter estimation. In *2016 IEEE International Symposium on Information Theory (ISIT)*, pages 1551–1555, 2016.
- [8] S. Axler. *Measure, integration and real analysis*, volume 282 of *Graduate Texts in Mathematics*. Springer, Cham, 2020.
- [9] D. Batenkov and N. Diab. Super-resolution of generalized spikes and spectra of confluent Vandermonde matrices. *Appl. Comput. Harmon. Anal.*, 65:181–208, 2023.
- [10] D. Batenkov, B. Diederichs, G. Goldman, and Y. Yomdin. The spectral properties of Vandermonde matrices with clustered nodes. *Linear Algebra Appl.*, 609:37–72, 2021.
- [11] D. Batenkov and G. Goldman. Single-exponential bounds for the smallest singular value of Vandermonde matrices in the sub-Rayleigh regime. *Appl. Comput. Harmon. Anal.*, 55:426–439, 2021.
- [12] D. Batenkov, G. Goldman, and Y. Yomdin. Super-resolution of near-colliding point sources. *Inf. Inference*, 10(2):515–572, 2021.
- [13] A. Beurling. Sur les intégrales de fourier absolument convergentes et leur application à une transformation fonctionnelle. In *9th Scand. Math. Congress*, pages 345–366, 1938.
- [14] A. Beurling. *The collected works of Arne Beurling. Vol. 2*. Contemporary Mathematicians. Birkhäuser Boston, Inc., Boston, MA, 1989. Harmonic analysis, Edited by L. Carleson, P. Malliavin, J. Neuberger and J. Wermer.

## Bibliography

- [15] Å. Björck. *Numerical methods in matrix computations*, volume 59 of *Texts in Applied Mathematics*. Springer, Cham, 2015.
- [16] M. Born and E. Wolf. *Principles of optics*. Pergamon Press, sixth edition, 1986. Reprinted (with corrections).
- [17] K. Bredies and H. K. Pikkarainen. Inverse problems in spaces of measures. *ESAIM Control Optim. Calc. Var.*, 19(1):190–218, 2013.
- [18] P. Breiding and N. Vannieuwenhoven. The condition number of Riemannian approximation problems. *SIAM J. Optim.*, 31(1):1049–1077, 2021.
- [19] P. Bürgisser and F. Cucker. *Condition*, volume 349 of *Grundlehren der mathematischen Wissenschaften*. Springer, Heidelberg, 2013.
- [20] P. L. Butzer and R. J. Nessel. *Fourier analysis and approximation*. Pure and Applied Mathematics, Vol. 40. Academic Press, New York-London, 1971.
- [21] C. A. Cabrelli and U. M. Molter. The Kantorovich metric for probability measures on the circle. *J. Comput. Appl. Math.*, 57(3):345–361, 1995.
- [22] E. Candès and C. Fernandez-Granda. Towards a mathematical theory of super-resolution. *Comm. Pure Appl. Math.*, 67(6):906–956, 2013.
- [23] J. Carruth, N. Elkies, F. Gonçalves, and M. Kelly. The Beurling-Selberg Box Minorant Problem via Linear Programming Bounds. *arXiv: Classical analysis*, 2022.
- [24] J. T. Carruth. *Extremal problems in Fourier analysis, Whitney’s theorem, and the interpolation of data*. PhD thesis, The University of Texas at Austin, 2019.
- [25] P. Catala, M. Hockmann, and S. Kunis. Sparse super resolution and its trigonometric approximation in the p-Wasserstein distance. *Proc. Appl. Math. Mech.*, 22(1), 2023.
- [26] P. Catala, M. Hockmann, S. Kunis, and M. Wageringel. Approximation and interpolation of singular measures by trigonometric polynomials. *arXiv: Numerical Analysis*, 2022.
- [27] S. Chen and A. Moitra. Algorithmic foundations for the diffraction limit. *ArXiv: Data Structures and Algorithms*, 2020.
- [28] S. Chen and A. Moitra. Algorithmic foundations for the diffraction limit. In *STOC ’21—Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 490–503. ACM, New York, 2021.
- [29] H. Cohn and N. Elkies. New upper bounds on sphere packings. I. *Ann. of Math. (2)*, 157(2):689–714, 2003.
- [30] A. Cuyt, W.-s. Lee, and X. Yang. On tensor decomposition, sparse interpolation and Padé approximation. *Jaen J. Approx.*, 8(1):33–58, 2016.
- [31] Y. de Castro and F. Gamboa. Exact Reconstruction using Beurling Minimal Extrapolation. *J. Math. Anal. Appl.*, 395(1):336–354, 2012.

- [32] Y. de Castro, F. Gamboa, D. Henrion, and J. Lasserre. Exact solutions to super resolution on semi-algebraic domains in higher dimensions. *IEEE Trans. Inform. Theory*, 63(1):621–630, 2017.
- [33] A. J. Den Dekker and A. Van den Bos. Resolution: a survey. *JOSA A*, 14(3):547–557, 1997.
- [34] Q. Denoyelle, V. Duval, and G. Peyré. Support recovery for sparse super-resolution of positive measures. *J. Fourier Anal. Appl.*, 23:1153–1194, 2017.
- [35] Q. Denoyelle, V. Duval, G. Peyré, and E. Soubies. The sliding Frank-Wolfe algorithm and its application to super-resolution microscopy. *Inverse Problems*, 36(1):014001, 42, 2020.
- [36] R. A. DeVore and G. G. Lorentz. *Constructive approximation*, volume 303 of *Grundlehren der mathematischen Wissenschaften*. Springer-Verlag, Berlin, 1993.
- [37] B. Diederichs. *Sparse Frequency Estimation : Stability and Algorithms*. PhD thesis, University of Hamburg, 2018.
- [38] B. Diederichs. Well-posedness of sparse frequency estimation. *arXiv: Numerical Analysis*, 2019.
- [39] D. L. Donoho. Superresolution via sparsity constraints. *SIAM J. Math. Anal.*, 23(5):1309–1331, 1992.
- [40] A. L. Dontchev and R. T. Rockafellar. *Implicit functions and solution mappings*. Springer Series in Operations Research and Financial Engineering. Springer, New York, second edition, 2014.
- [41] D. Dryanov and P. Petrov. Interpolation and  $L_1$ -approximation by trigonometric polynomials and blending functions. *J. Approx. Theory*, 164(8):1049–1064, 2012.
- [42] B. A. Dumitrescu. *Positive Trigonometric Polynomials and Signal Processing Applications*. Signals and Communication Technology. Springer International Publishing, 2017.
- [43] V. Duval and G. Peyré. Exact support recovery for sparse spikes deconvolution. *Found. Comput. Math.*, 15(5):1315–1355, 2015.
- [44] A. Eftekhari, T. Bendory, and G. Tang. Stable super-resolution of images: theoretical study. *Inf. Inference*, 10(1):161–193, 2021.
- [45] A. Eftekhari, J. Tanner, A. Thompson, B. Toader, and H. Tyagi. Sparse non-negative super-resolution—simplified and stabilised. *Appl. Comput. Harmon. Anal.*, 50:216–280, 2021.
- [46] M. Ehler, M. Gräf, S. Neumayer, and G. Steidl. Curve based approximation of measures on manifolds by discrepancy minimization. *Found. Comput. Math.*, 21(6):1595–1642, 2021.
- [47] M. Ehler, S. Kunis, T. Peter, and C. Richter. A randomized multivariate matrix pencil method for superresolution microscopy. *Elec. Trans. Numer. Anal.*, 51:63–74, 2019.

## Bibliography

- [48] D. Elbrächter, D. Perekrestenko, P. Grohs, and H. Bölcskei. Deep neural network approximation theory. *arXiv: Machine Learning*, 2019.
- [49] Z. Fan and J. Y. Li. Efficient algorithms for sparse moment problems without separation. *arXiv: Machine Learning*, 2022.
- [50] M. Fatemi, A. Amini, and M. Vetterli. Sampling and reconstruction of shapes with algebraic boundaries. *IEEE Trans. Signal Process.*, 64(22):5807–5818, 2016.
- [51] J. Favard. Sur les meilleurs procédés d’approximation de certaines classes de fonctions par des polynômes trigonométriques. *Bull. Sci. Math.*, 61:209–224, 1937.
- [52] C. Fernandez-Granda. Super-resolution of point sources via convex programming. *Inf. Inference*, 5(3):251–303, 2016.
- [53] M. Ferreira Da Costa and U. Mitra. On the Stability of Super-Resolution and a Beurling–Selberg Type Extremal Problem. In *2022 IEEE International Symposium on Information Theory (ISIT)*, pages 1737–1742, 2022.
- [54] M. Field. *Essential real analysis*. Springer Undergraduate Mathematics Series. Springer, Cham, 2017.
- [55] S. D. Fisher. Best approximation by polynomials. *J. Approx. Theory*, 21(1):43–59, 1977.
- [56] L. R. F.R.S. Investigations in optics, with special reference to the spectroscope. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 8(49):261–274, 1879.
- [57] W. Gautschi. On inverses of Vandermonde and confluent Vandermonde matrices. *Numer. Math.*, 4:117–123, 1962.
- [58] G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, third edition, 1996.
- [59] F. Gonçalves. A note on band-limited minorants of an Euclidean ball. *Proc. Amer. Math. Soc.*, 146(5):2063–2068, 2018.
- [60] D. V. Gorbachev. Extremal problems for entire functions of exponential spherical type. *Mat. Zametki*, 68(2):179–187, 2000.
- [61] L. Grafakos. *Classical Fourier analysis*, volume 249 of *Graduate Texts in Mathematics*. Springer, New York, third edition, 2014.
- [62] M. G. L. Gustafsson. Surpassing the lateral resolution limit by a factor of two using structured illumination microscopy. *J. Microsc.*, 198(2):82–87, 2000.
- [63] M. G. L. Gustafsson. Nonlinear structured-illumination microscopy: Wide-field fluorescence imaging with theoretically unlimited resolution. *Proc. Natl. Acad. Sci. U.S.A.*, 102(37):13081–13086, 2005.
- [64] H. He and D. Kressner. Randomized joint diagonalization of symmetric matrices. *arXiv: Numerical analysis*, 2022.

- [65] R. Heintzmann and C. G. Cremer. Laterally modulated excitation microscopy: improvement of resolution by using a diffraction grating. In I. J. Bigio, H. Schneck-enburger, J. Slavik, K. S. M.D., and P. M. Viallet, editors, *Optical Biopsies and Microscopic Techniques III*, volume 3568, pages 185 – 196. International Society for Optics and Photonics, SPIE, 1999.
- [66] R. Heintzmann, T. M. Jovin, and C. Cremer. Saturated patterned excitation microscopy—a concept for optical resolution improvement. *J. Opt. Soc. Am. A*, 19(8):1599–1609, 2002.
- [67] M. Hockmann and S. Kunis. Sparse super resolution is Lipschitz continuous. *arXiv: Numerical Analysis*, 2021.
- [68] M. Hockmann and S. Kunis. Short Communication: Weak Sparse Superresolution is Well-Conditioned. *SIAM J. Imaging Sci.*, 16(1):SC1–SC13, 2023.
- [69] M. Hockmann, S. Kunis, and R. Kurre. Towards a mathematical model for single molecule structured illumination microscopy. *Proc. Appl. Math. Mech.*, 20(1), 2021.
- [70] M. Hockmann, S. Kunis, and R. Kurre. Computational resolution in single molecule localization – impact of noise level and emitter density. *Biol. Chem.*, 404(5):427–431, 2023.
- [71] J. J. Holt and J. D. Vaaler. The Beurling-Selberg extremal functions for a ball in Euclidean space. *Duke Math. J.*, 83(1):202–248, 1996.
- [72] L. Hörmander. *The analysis of linear partial differential operators. I*. Springer Study Edition. Springer-Verlag, Berlin, second edition, 1990.
- [73] R. A. Horn and C. R. Johnson. *Matrix analysis*. Cambridge University Press, Cambridge, second edition, 2013.
- [74] Y. Hua and T. K. Sarkar. Matrix pencil method for estimating parameters of exponentially damped/undamped sinusoids in noise. *IEEE Trans. Acoust. Speech Signal Process.*, 38(5):814–824, 1990.
- [75] E. Ingerman, R. London, R. Heintzmann, and M. Gustafsson. Signal, noise and resolution in linear and nonlinear structured-illumination microscopy. *J. Microsc.*, 273(1):3–25, 2019.
- [76] A. E. Ingham. Some trigonometrical inequalities with applications to the theory of series. *Math. Z.*, 41(1):367–379, 1936.
- [77] D. Jackson. *The theory of approximation*, volume 11 of *American Mathematical Society Colloquium Publications*. American Mathematical Society, Providence, RI, 1994. Reprint of the 1930 original.
- [78] F. Johansson. Arb: efficient arbitrary-precision midpoint-radius interval arithmetic. *IEEE Trans. Comput.*, 66(8):1281–1292, 2017.
- [79] C. Jozs, J. Lasserre, and B. Mourrain. Sparse polynomial interpolation: Compressed sensing, super resolution, or Prony? *Adv. Comput. Math.*, 45(3):1401–1437, 2019.

## Bibliography

- [80] D. W. Kammler. *A first course in Fourier analysis*. Cambridge University Press, Cambridge, second edition, 2007.
- [81] H. Knirsch. *Optimal Hankel Structured Rank-1 Approximation*. PhD thesis, Georg-August-Universität Göttingen, 2022.
- [82] H. Knirsch, M. Petz, and G. Plonka. Optimal rank-1 Hankel approximation of matrices: Frobenius norm and spectral norm and Cadzow’s algorithm. *Linear Algebra Appl.*, 629:1–39, 2021.
- [83] V. Komornik and P. Loreti. *Fourier series in control theory*. Springer Monographs in Mathematics. Springer-Verlag, New York, 2005.
- [84] A. Kroó and D. Lubinsky. Christoffel functions and universality in the bulk for multivariate orthogonal polynomials. *Can. J. Math.*, 65(3):600–620, 2012.
- [85] A. Kroó and J. J. Swetits. On density of interpolation points, a Kadec-type theorem, and Saff’s principle of contamination in  $L_p$ -approximation. *Constr. Approx.*, 8(1):87–103, 1992.
- [86] S. Kunis, H. M. Möller, T. Peter, and U. von der Ohe. Prony’s method under an almost sharp multivariate Ingham inequality. *J. Fourier Anal. Appl.*, 24(5):1306–1318, 2018.
- [87] S. Kunis and D. Nagel. On the smallest singular value of multivariate Vandermonde matrices with clustered nodes. *Linear Algebra Appl.*, 604:1–20, 2020.
- [88] S. Kunis, D. Nagel, and A. Strotmann. Multivariate Vandermonde matrices with separated nodes on the unit circle are stable. *Appl. Comput. Harmon. Anal.*, 58:50–59, 2022.
- [89] S. Kunis, T. Peter, T. Römer, and U. von der Ohe. A multivariate generalization of Prony’s method. *Linear Algebra Appl.*, 490:31–47, 2016.
- [90] B. Kurmanbek and E. Robeva. Multivariate super-resolution without separation. *arXiv: Optimization and Control*, 2022.
- [91] R. Lakshmanan, A. Pichler, and D. Potts. Nonequispaced fast Fourier transform boost for the Sinkhorn algorithm. *Electron. Trans. Numer. Anal.*, 58:289–315, 2023.
- [92] A. Lal, C. Shan, and P. Xi. Structured illumination microscopy image reconstruction algorithm. *IEEE J. Sel. Top. Quantum Electron.*, 22:50–63, 2016.
- [93] J. Lasserre and E. Pauwels. The empirical Christoffel function with applications in data analysis. *Adv. Comput. Math.*, 45(3):1439–1468, 2019.
- [94] M. Laurent. *Sums of Squares, Moment Matrices and Optimization Over Polynomials*, pages 157–270. Springer New York, New York, NY, 2009.
- [95] M. Laurent and P. Rostalski. The approach of moments for polynomial equations. In *Handbook on semidefinite, conic and polynomial optimization*, volume 166 of *Internat. Ser. Oper. Res. Management Sci.*, pages 25–60. Springer, New York, 2012.
- [96] B. Laville, L. Blanc-Féraud, and G. Aubert. Off-the-grid variational sparse spike recovery: Methods and algorithms. *J. Imaging*, 7(12), 2021.



- [97] J. Lellmann, D. A. Lorenz, C. Schönlieb, and T. Valkonen. Imaging with Kantorovich-Rubinstein discrepancy. *SIAM J. Imaging Sci.*, 7(4):2833–2859, 2014.
- [98] N. Lev and J. Ortega-Cerdà. Equidistribution estimates for Fekete points on complex manifolds. *J. Eur. Math. Soc. (JEMS)*, 18(2):425–464, 2016.
- [99] W. Li, W. Liao, and A. Fannjiang. Super-resolution limit of the ESPRIT algorithm. *IEEE Trans. Inform. Theory*, 66(7):4593–4608, 2020.
- [100] W. Li, Z. Zhu, W. Gao, and W. Liao. Stability and super-resolution of MUSIC and ESPRIT for multi-snapshot spectral estimation. *IEEE Trans. Signal Process.*, 70:4555–4570, 2022.
- [101] W. Liao and A. Fannjiang. MUSIC for single-snapshot spectral estimation: stability and super-resolution. *Appl. Comput. Harmon. Anal.*, 40:33–67, 2016.
- [102] F. Littmann. Quadrature and extremal bandlimited functions. *SIAM J. Math. Anal.*, 45(2):732–747, 2013.
- [103] P. Liu. *Mathematical Theory of Computational Resolution Limit and Efficient Fast Algorithms for Super-resolution*. PhD thesis, Hong Kong University of Science and Technology, 2021.
- [104] P. Liu and H. Ammari. Nearly optimal resolution estimate for the two-dimensional super-resolution and a new algorithm for direction of arrival estimation with uniform rectangular array. *arXiv: Image and Video Processing*, 2022.
- [105] P. Liu and H. Ammari. Super-resolution of positive near-colliding point sources. *arXiv: Image and Video Processing*, 2022.
- [106] P. Liu, Y. He, and H. Ammari. A mathematical theory of resolution limits for super-resolution of positive sources. *arXiv: Image and Video Processing*, 2022.
- [107] P. Liu, S. Yu, O. Sabet, L. Pelkmans, and H. Ammari. Mathematical foundation of sparsity-based multi-illumination super-resolution. *arXiv: Image and Video Processing*, 2022.
- [108] P. Liu and H. Zhang. A mathematical theory of computational resolution limit in multi-dimensional spaces. *Inverse Problems*, 37(10):Paper No. 104001, 30, 2021.
- [109] P. Liu and H. Zhang. A theory of computational resolution limit for line spectral estimation. *IEEE Trans. Inf. Theory*, 67(7):4812–4827, 2021.
- [110] D. Manolakis, V. Ingle, and S. Kogon. *Statistical and Adaptive Signal Processing*. ARTECH, 2005.
- [111] S. Marx, E. Pauwels, T. Weisser, D. Henrion, and J. Lasserre. Semi-algebraic approximation using Christoffel-darboux kernel. *Constr. Approx.*, 54:391–429, 2021.
- [112] H. Mehta. The  $L^1$  norms of de la Vallée Poussin kernels. *J. Math. Anal. Appl.*, 422(2):825–837, 2015.
- [113] H. N. Mhaskar. Super-resolution meets machine learning: approximation of measures. *J. Fourier Anal. Appl.*, 25(6):3104–3122, 2019.

## Bibliography

- [114] B. S. Mityagin. The zero set of a real analytic function. *Mat. Zametki*, 107(3):473–475, 2020.
- [115] L. Möckl, D. C. Lamb, and C. Bräuchle. Super-resolved fluorescence microscopy: Nobel Prize in Chemistry 2014 for Eric Betzig, Stefan Hell, and William E. Moerner. *Angew. Chem., Int. Ed.*, 53(51):13972 – 13977, 2014.
- [116] W. Moerner and D. Fromm. Methods of single-molecule fluorescence spectroscopy and microscopy. *Rev. Sci. Instr.*, 74(8):3597–3619, 2003.
- [117] A. Moitra. Super-resolution, extremal functions and the condition number of Vandermonde matrices. In *STOC’15—Proceedings of the 2015 ACM Symposium on Theory of Computing*, pages 821–830. ACM, New York, 2015.
- [118] E. Moskona, P. Petrushev, and E. B. Saff. The Gibbs phenomenon for best  $L_1$ -trigonometric polynomial approximation. *Constr. Approx.*, 11(3):391–416, 1995.
- [119] B. Mourrain. Polynomial-exponential decomposition from moments. *Found. Comput. Math.*, 18(6):1435–1492, 2018.
- [120] D. Nagel. *The condition number of Vandermonde matrices and its application to the stability analysis of a subspace method*. PhD thesis, Osnabrueck University, 2020.
- [121] E. Nehme, D. Freedman, R. Gordon, B. Ferdman, L. E. Weiss, O. Alalouf, T. Naor, R. Orange, T. Michaeli, and Y. Shechtman. DeepSTORM3D: dense 3D localization microscopy and PSF design by deep learning. *Nature methods*, 17(7):734–740, 2020.
- [122] E. Nehme, L. E. Weiss, T. Michaeli, and Y. Shechtman. Deep-STORM: super-resolution single-molecule microscopy by deep learning. *Optica*, 5(4):458–464, Apr 2018.
- [123] P. Nevai. Géza Freud, orthogonal polynomials and Christoffel functions. a case study. *J. Approx. Theory*, 48(1):3–167, 1986.
- [124] G. Ongie and M. Jacob. Off-the-grid recovery of piecewise constant images from few Fourier samples. *SIAM J. Imaging Sci.*, 9(3):1004–1041, 2016.
- [125] M. Ovesný, P. Křížek, J. Borkovec, Z. Svindrych, and G. M. Hagen. ThunderSTORM: a comprehensive ImageJ plug-in for PALM and STORM data analysis and super-resolution imaging. *Bioinformatics*, 30(16):2389–2390, 2014.
- [126] P. Pakrooh, A. Pezeshki, L. L. Scharf, D. Cochran, and S. D. Howard. Analysis of Fisher information and the Cramér-Rao bound for nonlinear parameter estimation after random compression. *IEEE Trans. Signal Process.*, 63(23):6423–6428, 2015.
- [127] H. Pan, T. Blu, and P. Dragotti. Sampling curves with finite rate of innovation. *IEEE Trans. Signal Process.*, 62(2):458–471, 2014.
- [128] E. Pauwels, M. Putinar, and J.-B. Lasserre. Data analysis from empirical moments and the Christoffel function. *Found. Comput. Math.*, 21(1):243–273, 2021.
- [129] G. Peyré and M. Cuturi. Computational optimal transport: With applications to data science. *Found. Trends Mach. Learn.*, 11(5-6):355–607, 2019.

- [130] B. Piccoli and F. Rossi. Generalized Wasserstein distance and its application to transport equations with source. *Arch. Ration. Mech. Anal.*, 211(1):335–358, 2014.
- [131] B. Piccoli, F. Rossi, and M. Tournus. A Wasserstein norm for signed measures, with application to nonlocal transport equation with source term. *arXiv: Analysis of PDE*, 2019.
- [132] G. Plonka, D. Potts, G. Steidl, and M. Tasche. *Numerical Fourier analysis*. Applied and Numerical Harmonic Analysis. Birkhäuser/Springer, Cham, 2018.
- [133] G. Plonka and M. Tasche. Prony methods for recovery of structured functions. *GAMM-Mitt.*, 37(2):239–258, 2014.
- [134] C. Poon and G. Peyré. Multi-dimensional sparse super-resolution. *SIAM J. Math. Anal.*, 51(1):1–44, 2018.
- [135] D. Potts and M. Tasche. Error estimates for the ESPRIT algorithm. In *Large truncated Toeplitz matrices, Toeplitz operators, and related topics*, volume 259 of *Oper. Theory Adv. Appl.*, pages 621–648. Birkhäuser/Springer, Cham, 2017.
- [136] R. Prony. Essai experimentable et analytique: Sur les lois de la Dilatabilité des fluides élastiques et sur celles de la Force expansive de la vapeur de l’eau et de la vapeur de l’alkool, à différentes températures. *Journal de l’École Polytechnique Floréal et Plairial*, 1:24–76, 1795.
- [137] R. Roy and T. Kailath. ESPRIT-estimation of signal parameters via rotational invariance techniques. *IEEE Trans. Acoustics Speech Signal Process.*, 37(7):984–995, 1989.
- [138] W. Rudin. *Real and complex analysis*. McGraw-Hill Book Co., New York, third edition, 1987.
- [139] M. J. Rust, M. Bates, and X. Zhuang. Sub-diffraction-limit imaging by stochastic optical reconstruction microscopy (STORM). *Nature Methods*, 3:793–796, 2006.
- [140] S. Sahnoun, K. Usevich, and P. Comon. Multidimensional ESPRIT for damped and undamped signals: Algorithm, computations, and perturbation analysis. *IEEE Trans. Signal Proc.*, 65(22):5897–5910, 2017.
- [141] F. Santambrogio. *Optimal transport for applied mathematicians*, volume 87 of *Progress in Nonlinear Differential Equations and their Applications*. Birkhäuser/Springer, Cham, 2015.
- [142] T. Sauer. Prony’s method in several variables. *Numer. Math.*, 136:411–438, 2017.
- [143] M. Schmidt. minFunc: unconstrained differentiable multivariate optimization in Matlab. <http://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>, 2005.
- [144] R. Schmidt. Multiple emitter location and signal parameter estimation. *IEEE Trans. Antennas Propagation*, 34(3):276–280, 1986.
- [145] A. Schuster. *An introduction to the theory of optics*. Edward Arnold, 1904.
- [146] A. Selberg. *Collected papers. I*. Springer Collected Works in Mathematics. Springer, Heidelberg, 2014.

## Bibliography

- [147] C. M. Sparrow. On spectroscopic resolving power. *Astrophysical Journal*, vol. 44, p. 76, 44:76, 1916.
- [148] A. Speiser, L.-R. Müller, U. Matti, C. J. Obara, W. R. Legant, A. Kreshuk, J. H. Macke, J. Ries, and S. C. Turaga. Deep learning enables fast and dense single-molecule localization with high accuracy. *Nat. Methods.*, 18(9):1082–1090, 2021.
- [149] L. N. Trefethen and D. Bau, III. *Numerical linear algebra*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997.
- [150] J. D. Vaaler. Some extremal functions in Fourier analysis. *Bull. Amer. Math. Soc. (N.S.)*, 12(2):183–216, 1985.
- [151] C. Villani. *Optimal transport - Old and new*, volume 338 of *Grundlehren der Mathematischen Wissenschaften*. Springer-Verlag, Berlin, 2009.
- [152] M. Wageringel. Truncated moment problems on positive-dimensional algebraic varieties. *arXiv: Commutative Algebra*, 2022.
- [153] J. Wallis, T. Miller, C. Lerner, and E. Kleerup. Three-dimensional display in nuclear medicine. *IEEE Transactions on Medical Imaging*, 8(4):297–230, 1989.
- [154] G. N. Watson. *A Treatise on the Theory of Bessel Functions*. Cambridge University Press, Cambridge, England; The Macmillan Company, New York, 1944.
- [155] P.-A. Wedin. Perturbation bounds in connection with singular value decomposition. *Nordisk Tidskr. Informationsbehandling (BIT)*, 12:99–111, 1972.
- [156] P.-A. Wedin. Perturbation theory for pseudo-inverses. *Nordisk Tidskr. Informationsbehandling (BIT)*, 13:217–232, 1973.
- [157] D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.
- [158] D.-X. Zhou. Universality of deep convolutional neural networks. *Appl. Comput. Harmon. Anal.*, 48(2):787–794, 2020.

# List of Figures

1	Diffraction limits in the literature . . . . .	6
1.1	Bessel functions . . . . .	21
1.2	Extremal functions . . . . .	23
2.1	Separation distance . . . . .	32
2.2	Manifold of moments for sparse measures . . . . .	32
2.3	Admissible functions . . . . .	40
2.4	Box minorant function . . . . .	43
2.5	Decomposition of the node set . . . . .	45
2.6	Set-valued reconstruction map . . . . .	51
2.7	Alternating node sets in lower bound . . . . .	54
2.8	Decomposition of hexagonal lattice node set . . . . .	56
2.9	Bivariate diffraction limit . . . . .	62
2.10	Pair clustering configuration . . . . .	70
3.1	Example measure 1D . . . . .	76
3.2	Example measure in 2D . . . . .	76
3.3	Best approximation on univariate torus . . . . .	87
3.4	Bounds on signal polynomial . . . . .	95
3.5	Rates for polynomials in a numerical example . . . . .	98
3.6	Christoffel function without noise . . . . .	102
3.7	Parameter choice for Christoffel function . . . . .	112
4.1	Principle of STORM . . . . .	116
4.2	STORM results . . . . .	121
4.3	Linear 2D-SIM . . . . .	124
4.4	Filter function SIM . . . . .	127
4.5	Results of SIM algorithms . . . . .	128

Except the logo of the University of Osnabrück, all figures have been created with MATLAB and the `TikZ` package. The code to reproduce the numerical experiments and the figures created in MATLAB is publicly available under <https://github.com/MHockmann/Dissertation.git>.

## List of Algorithms

1	Support approximation by Christoffel function . . . . .	113
2	STORM algorithm with signal polynomial . . . . .	122
3	SIM with separation and recombination in spatial domain. . . . .	127
4	SIM with separation in spatial and recombination in Fourier domain . . . .	128

## Glossary of symbols

$\mathbb{N}$	natural numbers $\{0, 1, 2, \dots\}$
$\mathbb{Z}$	integers
$\mathbb{R}$	real numbers
$\mathbb{C}$	complex numbers
$\mathbb{T}$	torus
$B_q(0)$	ball of radius $q$ around zero
$\Re(z)$	real part of complex number $z$
$\Im(z)$	imaginary part of complex number $z$
$\bar{z}$	complex conjugate of complex number $z$
$ Y $	cardinality of set $Y$
$\alpha^d$	$d$ -fold repetition of vector $\alpha$
$v^\top, v^*$	transpose and conjugate transpose
$\langle \cdot, \cdot \rangle$	inner product
$\lambda_j(A)$	$j$ th eigenvalue of $A$
$\lambda_{\min}(A), \lambda_{\max}(A)$	smallest and largest eigenvalue of Hermitian matrix $A$
$\sigma_j(A)$	$j$ th singular value of $A$
$\sigma_{\min}(A), \sigma_{\max}(A)$	smallest and largest singular value of $A$
$\text{img } A$	range of $A$
$\ker A$	kernel of $A$
$\text{diag}(v)$	diagonal matrix with diagonal given by vector $v$
$A^\dagger$	pseudo inverse of $A$
$\text{Tr}(A)$	trace of a matrix $A$
$\text{dist}$	distance function
$\text{proj}_U(x)$	orthogonal projection of $x$ onto subspace $U$
$\ \cdot\ _2$	Euclidean norm of vector or matrix
$\ \cdot\ _\infty, \ \cdot\ _p$	maximum norm or $p$ -norm of vector or matrix

$\ \cdot\ _F$	Frobenius norm of matrix
$\kappa_{\text{abs}}$	absolute condition number
$\tilde{\kappa}_{\text{abs}}$	condition number of recovery from image data
$\mathcal{A}$	Vandermonde matrix
$\tilde{\mathcal{A}}_s$	confluent Vandermonde matrix in dimension $s$
$\mathbf{T}_n$	moment matrix
$\mathbf{T}_{0,n}$	ground truth moment matrix
$\mathbf{T}_{\rho,n}$	noise moment matrix
$\mathbf{T}_{\tilde{r},n}$	truncated moment matrix
$\mathbb{A}_n$	Vandermonde matrix up to Euclidean degree $n/2$
$\mathbf{U}_n, \mathbf{\Sigma}_n, \mathbf{V}_n$	matrices from SVD of $\mathbf{T}_n$
$W$	matrix containing weights
$L^1(\mathbb{T}^d), L^2(\mathbb{T}^d)$	space of integrable or square integrable functions
$\mathcal{B}_n(\mathbb{R}^d)$	space of bandlimited functions
$C(\mathbb{T}^d)$	space of continuous functions
$C_c(\mathbb{T}^d)$	space of continuous functions with compact support
$L^\infty(\mathbb{T}^d)$	space of a.e. bounded functions
$\hat{f}$	Fourier transform of $f$
$f * g$	convolution of $f$ and $g$
$\text{Lip}(f)$	Lipschitz constant of $f$
$\partial$	partial derivative
$\Delta$	Laplace operator
$\text{sgn}$	sign function
$h$	Point-spread-function
$I$	illumination function
$J_\nu$	Bessel function of the first kind
$j_{\nu,k}$	$k$ th smallest positive zero of $J_\nu$
$\psi, \psi_\tau, \psi_{\tau,n}$	admissible function
$\mathbb{1}_C$	indicator function of set $C$
$\ \mu\ _{\text{TV}}$	total variation of $\mu$
$\mathcal{M}(\mathcal{X})$	space of complex Borel measures with finite total variation
$\mathcal{M}_{\mathbb{R}}(\mathcal{X})$	space of signed measures
$\mathcal{M}_+(\mathcal{X})$	space of nonnegative measures
$\mathcal{M}_{+,1}(\mathcal{X})$	space of probability measures
$\mu_k \rightharpoonup \mu$	weak convergence
$T_{\#}\mu$	push-forward measure
$\Pi(\mu, \nu)$	space of coupling of $\mu$ and $\nu$
$W_p$	$p$ -Wasserstein distance
$\text{supp}$	support of measure or function
$\alpha_t$	weights of discrete measure
$t$	node of discrete measure
$\delta_t$	Dirac measure at $t$
$Y$	support of discrete measure
$\rho, \hat{\rho}$	noise in spatial or Fourier domain
$\varrho$	noise level
$\ \cdot\ _{\mathbb{T}^d}$	wrap around distance

Glossary of symbols

$q, \text{sep } Y$	wrap around separation
$\text{clusep } Y$	cluster separation
$\mathcal{M}(q)$	$q$ -separated complex measures
$\hat{\mathcal{M}}^n(q)$	truncated moment set
$M$	upper bound on number of nodes
$\alpha_{\min}, \alpha_{\max}$	smallest and largest weight in absolute value
$Y_1, Y_2, Y_3$	decomposition of node sets from Theorem 2.2.8
$\mathbb{E}$	expected value
$J(\theta)$	Fisher information matrix
$\mathcal{CN}$	complex normal distribution
$S$	number of frames
$n$	bandlimit or parameter for truncation of moments
$\mathcal{I}$	sampling set $\mathcal{I} = \{k \in \mathbb{Z}^d : \ k\ _2 \leq n\}$
$\mathcal{R}$	reconstruction map
$\tilde{\mathcal{R}}$	image data reconstruction map
$\Omega_d$	optimal transition constant
$J$	sampling parameter
$\Delta$	field-of-view parameter
$n' = \gamma n$	truncation parameter in Fourier domain
$\tilde{\Omega}_2$	optimal transition constant for image recovery problem
$\tilde{\Omega}_{\text{SIM},2}$	optimal transition constant for SIM
$\mathcal{P}^{n,d,\infty}$	space of polynomials with maximal degree $n$
$\mathcal{P}^{n,d,2}$	space of polynomials with Euclidean degree $n$
$N$	dimension of $\mathcal{P}^{n/2,d,2}$
$D_n$	Dirichlet kernel
$d_n$	modified Dirichlet kernel
$F_n$	Fejér kernel
$J_n$	Jackson kernel
$D_{\text{rad},n}(x)$	radial Dirichlet kernel
$F_{\text{rad},n}(x)$	radial Fejér kernel
$p_n$	convolution of $\mu$ with $F_n$
$K$	reproducing kernel
$(H, \mathcal{D})$	reproducing kernel Hilbert space with discrepancy $\mathcal{D}$
$\mathcal{B}_1$	Bernoulli spline of degree 1
$e_x^{(n)}$	vector of monomials
$\mathbf{p}$	vector of coefficients of general polynomial $p$
$V(\ker \mathbf{T}_n)$	variety spanned by polynomials in $\ker \mathbf{T}_n$
$p_{1,n}$	signal polynomial
$\tilde{p}_{1,n}$	perturbed signal polynomial
$p_{0,n}$	noise polynomial
$u_j^{(n)}, v_j^{(n)}$	polynomials with coefficients $\mathbf{u}_j^{(n)}$ or $\mathbf{v}_j^{(n)}$ respectively
$q_\mu$	Christoffel function
$q_{\varepsilon,n}$	regularised Christoffel function
$\tilde{q}_{\varepsilon,n}$	perturbed, regularised Christoffel function