



**HAL**  
open science

# Spurious minimizers in non uniform Fourier sampling optimization

Alban Gossard, Pierre Weiss, Frédéric de Gournay

► **To cite this version:**

Alban Gossard, Pierre Weiss, Frédéric de Gournay. Spurious minimizers in non uniform Fourier sampling optimization. 2021. hal-03212145v1

**HAL Id: hal-03212145**

**<https://hal.science/hal-03212145v1>**

Preprint submitted on 29 Apr 2021 (v1), last revised 20 Jul 2022 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Spurious minimizers in non uniform Fourier sampling optimization

Alban Gossard  
alban.paul.gossard@gmail.com

Frédéric de Gournay  
degourna@insa-toulouse.fr

Pierre Weiss  
pierre.armand.weiss@gmail.com

Institut de Mathématiques de Toulouse

April 29, 2021

**Abstract**— A recent trend in the signal/image processing literature is the optimization of Fourier sampling schemes for specific datasets of signals. In this paper, we explain why choosing optimal non Cartesian Fourier sampling patterns is a difficult nonconvex problem by bringing to light two optimization issues. The first one is the existence of a combinatorial number of spurious minimizers for a generic class of signals. The second one is a vanishing gradient effect for the high frequencies. We conclude the paper by showing how using large datasets can mitigate first effect and illustrate experimentally the benefits of using stochastic gradient algorithms with a variable metric.

## 1 Introduction

Finding efficient Fourier sampling schemes is a critical issue in communications and imaging. This led to various theories including the celebrated Shannon-Nyquist theorems for bandlimited signals and compressed sensing for sparse signals. Unfortunately - in most practical cases - the signals to reconstruct are quite loosely described by these generic classes. For instance, magnetic resonance images of brains or knees have a rich structure due to the physiological underlying object. It is therefore tempting to optimize a sampling scheme directly for a given dataset rather than relying on a rough mathematical model. The recent progresses in GPU programming, automatic differentiation and machine learning make this idea even more tantalizing. In the sole field of MRI, the following list of references [7, 8, 11, 2, 15, 12, 13, 12, 6] illustrates this novel trend.

Unfortunately, most of the above works report (more or less explicitly) optimization issues. Fig. 1 illustrates one of them. In this example, we tried to optimize a sampling scheme for a single image from the fastMRI challenge [14]. To this end, we minimize the  $\ell^2$  reconstruction error using a sim-

ple back-projection reconstructor with a subsampling factor of 2. The trajectory of a gradient descent is displayed in Fig. 1b. As can be seen, the final sampling set covers approximately uniformly the Fourier domain, while we would expect the low frequencies to be sampled more densely. This likely highlights the presence of a spurious minimizer.

The aim of this paper is to explain this phenomenon from a mathematical perspective and to bring some solutions to mitigate the difficulties. We focus on optimization schemes that continuously optimize the positions of some sampling locations. These techniques have the advantage of not relying on a grid, which is an essential feature for various applications such as magnetic resonance imaging or radio-interferometry. In addition, they spark the hope of avoiding the curse of dimensionality encountered in combinatorial problems. We show that this dream is not realistic, but that the situation improves by considering large signals datasets and specific variable metric techniques. We conclude the paper by illustrating our findings on 1D experiments.

## 2 Notation

In this paper, we will focus on discrete 1D signals, for the ease of exposition. However, the main arguments apply to arbitrary dimensions and continuous signals as well.

We consider a signal  $u$  as a vector of  $\mathbb{C}^N$  with  $N \in 2\mathbb{N}$ . We let  $\mathcal{N} = \llbracket -\frac{N}{2}, \frac{N}{2} - 1 \rrbracket$ . An alternative way to represent a signal  $u \in \mathbb{C}^N$  is to use a discrete measure  $\mu$  of the form:

$$\mu = \sum_{n \in \mathcal{N}} u_n \delta_{\frac{n}{N}}. \quad (1)$$

Given a location  $\xi \in \mathbb{R}$ , we define:

$$\hat{u}(\xi) \stackrel{\text{def}}{=} \frac{1}{\sqrt{N}} \sum_{n \in \mathcal{N}} u_n e^{-2i\pi \langle \xi, \frac{n}{N} \rangle}, \quad (2)$$

which can be seen as the continuous Fourier transform of the measure  $\mu$ . We consider  $\Xi =$

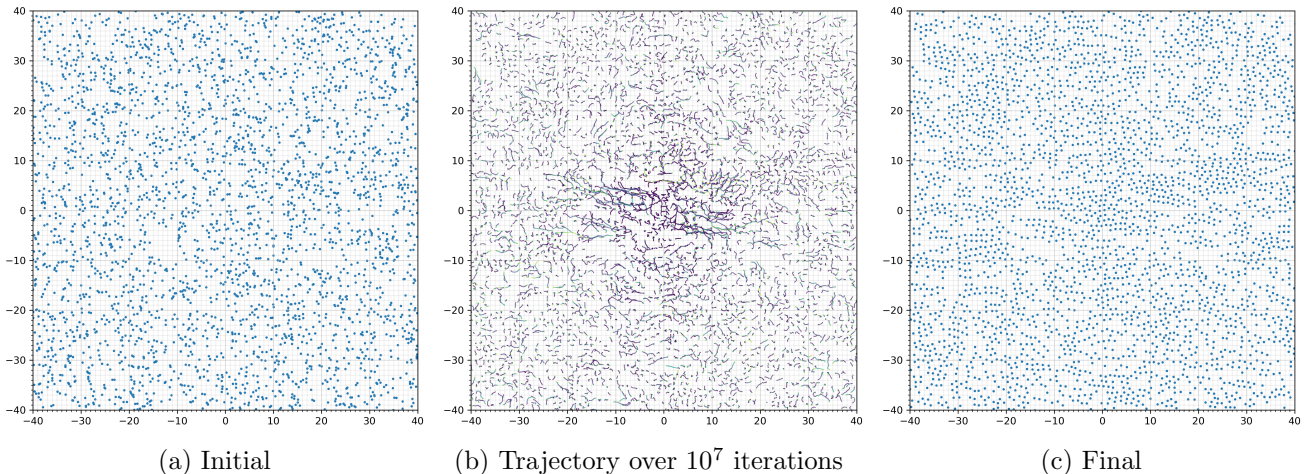


Figure 1: A typical sampling optimization trajectory. Starting from the sampling configuration on the left (uniform point process), we obtain the sampling scheme on the right after  $10^7$  iterations. The trajectory in the center corresponds to the  $10^7$  iterations of a gradient descent with fixed step size. Notice that the points clusters have disappeared, but that the scheme is still essentially uniform, while we would expect the low frequencies to be sampled more densely.

$[\xi_1, \dots, \xi_M] \in \mathbb{R}^M$  a set of  $M$  locations. The Fourier transform  $\hat{u}(\Xi) \in \mathbb{C}^M$  at the locations  $\Xi$  can be written as a matrix-vector product of the form  $\hat{u}(\Xi) = A(\Xi)^* u$  with the normalized Vandermonde matrix  $A(\Xi) \in \mathbb{C}^{N \times M}$  defined by

$$A(\Xi)_{n,m} \stackrel{\text{def}}{=} \frac{1}{\sqrt{N}} e^{2i\pi \langle \xi_m, \frac{n}{N} \rangle}.$$

In what follows, we let  $a(\xi) \in \mathbb{C}^N$  denote the vector defined for all  $n \in \mathcal{N}$  by

$$a(\xi)[n] \stackrel{\text{def}}{=} \frac{1}{\sqrt{N}} e^{2i\pi \langle \xi, \frac{n}{N} \rangle},$$

so that

$$A(\Xi) = [a(\xi_1), \dots, a(\xi_M)].$$

The matrix  $A(\Xi)^*$  can be seen as the nonuniform Fourier transform [10] from the grid to the set of sampling locations  $\Xi$ . We let  $(A(\Xi)^*)^+$  denote the pseudo-inverse of  $A(\Xi)^*$ .

### 3 Preliminaries

Below, we first describe the precise mathematical setting and then turn to some preliminary results.

#### 3.1 The setting

Let  $u \in \mathbb{C}^N$  denote a signal. We assume that a sampling device allows to pick  $M$  frequencies  $\xi_1, \dots, \xi_M$  in  $\mathbb{R}$ , yielding the set of measurements  $y = A(\Xi)^* u$ . A vast amount of reconstruction techniques have been designed in the literature to reconstruct  $u$  from  $y$ . A generic reconstructor can be defined as

a mapping  $R : (\mathbb{C}^M \times \mathbb{R}^M) \rightarrow \mathbb{C}^N$  that takes as an input a measurement  $y \in \mathbb{C}^M$  and a sampling scheme  $\Xi \in \mathbb{R}^M$  and outputs a reconstructed signal  $R(y, \Xi)$ . Given a collection of signals  $u_1, \dots, u_P$  and a reconstructor  $R$ , a natural framework to find the best sampling scheme  $\Xi$  is to solve the following optimization problem:

$$\inf_{\Xi \in \mathbb{R}^M} \frac{1}{2P} \sum_{p=1}^P \|R(A(\Xi)^* u_p, \Xi) - u_p\|_2^2. \quad (3)$$

This problem can be attacked with first order methods that continuously optimize the sampling locations  $\xi_m$ , see for instance [13, 6, 12]. In this work, we will concentrate on two simple linear reconstruction methods:

**The back-projection method** which consists in defining the reconstructor as

$$R(y, \Xi) = A(\Xi)y. \quad (4)$$

**The pseudo-inverse method** which consists in defining the reconstructor as

$$R(y, \Xi) = (A(\Xi)^*)^+ y. \quad (5)$$

Both techniques are quite popular in the literature. In fact, they coincide whenever  $\Xi$  is a subgrid (see definition hereafter) since in that case, the Fourier atoms are pairwise orthogonal. We restrict our analysis to linear reconstructors of the type (4) and (5) for simplicity reasons. Numerical experiments reveal that the optimization issues raised in Theorems 1 and 2 also apply to nonlinear reconstructors

such as sparsity promoting convex penalties. However, the techniques used in the proofs do not seem to easily extend to this framework.

We first analyze the problem with a single image  $u$  in the dataset, i.e.  $P = 1$ . Let us define the following two cost functions.

$$J_1(\Xi) \stackrel{\text{def}}{=} \frac{1}{2} \|A(\Xi)A(\Xi)^*u - u\|_2^2 \quad (6)$$

and

$$J_2(\Xi) \stackrel{\text{def}}{=} \frac{1}{2} \|(A(\Xi)^*)^+ A(\Xi)^*u - u\|_2^2 \quad (7)$$

Minimizing  $J_1$  allows to optimize the sampling scheme associated to the back-projection while minimizing  $J_2$  allows to optimize the sampling scheme for the pseudo-inverse.

### 3.2 Elementary observations

We will make use of the following definitions.

**Definition 1** (The min distance). *Given a set of sampling points  $\Xi$ , the min distance  $\text{md}(\Xi)$  is defined by*

$$\text{md}(\Xi) \stackrel{\text{def}}{=} \min_{m \neq m'} \text{dist}(\xi_m, \xi_{m'})$$

where  $\text{dist}$  is the distance on the torus defined for  $(\xi_1, \xi_2) \in \mathbb{R}^2$  as

$$\text{dist}(\xi_1, \xi_2) \stackrel{\text{def}}{=} \inf_{k \in \mathbb{Z}} \|\xi_1 - \xi_2 - kN\|_\infty. \quad (8)$$

**Definition 2** (Subgrid). *Throughout the paper, we say that  $\Xi \in [-N/2, N/2]^M$  is a subgrid if  $\xi_m - \xi_{m'} \in \mathbb{Z}^*$  for all  $m \neq m'$ .*

In what follows  $J$  can denote either  $J_1$  or  $J_2$  defined in (6) and (14).

**Proposition 1** ( $J$  is  $N$ -periodic). *We have*

$$J(\Xi \bmod N) = J(\Xi). \quad (9)$$

*Proof.* Let  $n = kN$  with  $k \in \mathbb{N}$ . The proof simply stems from the fact that  $a(\xi + n) = a(\xi)$ .  $\square$

The previous proposition shows that we can restrict our attention to frequencies  $\xi$  belonging to the set  $[-N/2, N/2[$ .

**Proposition 2** (Existence of minimizers). *For any  $M \in \mathbb{N}$  and any  $u \in \mathbb{C}^N$ , there exists at least one minimizer of  $J$  on  $[-N/2, N/2[$ .*

*Proof.* We start by noticing that  $J$  is a  $C^\infty$  function since it is defined as a composition of  $C^\infty$  functions. Hence it is also continuous on  $[-N/2, N/2[$ . This yields the existence of at least one minimizer.  $\square$

Now we proceed to a reformulation of the problem by rearranging the terms involved in the definition of  $J$  defined in (6). To that end, let us introduce the following function

$$F(\Xi) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{m=1}^M |\hat{u}(\xi_m)|^2 = \frac{1}{2} \|\hat{u}(\Xi)\|_2^2, \quad (10)$$

which somehow measures the energy captured within a sampling scheme  $\Xi$ . We also introduce the functions

$$G_1(\Xi) \stackrel{\text{def}}{=} \frac{1}{2} \langle (A(\Xi)^*A(\Xi) - \text{Id}) \hat{u}(\Xi), \hat{u}(\Xi) \rangle, \quad (11)$$

and

$$G_2(\Xi) \stackrel{\text{def}}{=} \frac{1}{2} \langle (\text{Id} - (A(\Xi)^*A(\Xi))^+) \hat{u}(\Xi), \hat{u}(\Xi) \rangle. \quad (12)$$

**Proposition 3.** *We have*

$$J_1(\Xi) = \frac{1}{2} \|u\|_2^2 - F(\Xi) + G_1(\Xi) \quad (13)$$

and

$$J_2(\Xi) = \frac{1}{2} \|u\|_2^2 - F(\Xi) + G_2(\Xi). \quad (14)$$

In particular, when  $\Xi$  is a subgrid, we have

$$J_1(\Xi) = J_2(\Xi) = \frac{1}{2} \|u\|_2^2 - \frac{1}{2} \|\hat{u}(\Xi)\|_2^2.$$

*Proof.* The proof of (13) and (14) is postponed to Section 7.1. When  $\Xi$  is a subgrid, the matrix  $A(\Xi)^*A(\Xi) - \text{Id}$  vanishes and so do  $G_1$  and  $G_2$ .  $\square$

## 4 Theoretical issues

In this section, we give the main theoretical results of the paper.

### 4.1 Spurious minimizers

The aim of this section is to illustrate common situations where the functions  $J_1$  and  $J_2$  both possess a combinatorial number of minimizers. Following Proposition 3, we construct examples where the function  $F$  is very oscillatory, while  $G_1$  and  $G_2$  are of small amplitude.

**Theorem 1** (A combinatorial number of minimizers). *Set a number of samples  $M \in \mathbb{N}$  and consider a vector  $u \in \mathbb{C}^N$  such that the following properties are verified*

1. *The modulus  $|\hat{u}|^2$  possesses a subset of  $K \geq M$  strict maximizers  $Z = \{\zeta_1, \dots, \zeta_K\}$  separated by a distance at least  $\delta = \text{md}(Z)$ .*

2. The modulus  $|\hat{u}|^2$  is locally strictly concave for each  $\zeta_k$ :

$$(|\hat{u}|^2)''(\xi) \leq -c, \quad \forall \xi \in [\zeta_k - r, \zeta_k + r]$$

for some radius  $0 < r < \delta/2$  and constant  $c > 0$ .

3. For any set  $\Xi$  of  $M$  distinct points in  $Z$ ,  $F(\Xi) \leq C$  for some constant  $C \geq 0$ .

Then under the condition

$$C < \frac{cr^2(\delta - 2r)}{4}, \quad (15)$$

the function  $J_1$  possess at least  $\binom{M}{K} \cdot M!$  local maximizers. The same result holds for  $J_2$  under the conditions  $\delta > 1$  and  $C < \frac{cr^2(\delta - 1 - 2r)}{4}$ .

The proof of Theorem 1 is postponed to Section 7.2. The conditions in Theorem 1 may look cryptic at first sight. Let us show a simple example of a function  $u$  that verifies the hypotheses and leads to a huge number of critical points.

**Corollary 1.** Assume that  $N \in 4\mathbb{N}$  and define  $u \in \mathbb{C}^N$  as follows

$$u[n] = \begin{cases} \sqrt{N}/2 & \text{if } n = \pm N/4, \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

Let  $M = \lfloor 0.11\sqrt{N} \rfloor$ , then  $J_1$  possesses a number of minimizers, which is equivalent to  $\alpha e^{\beta\sqrt{N}} N^{-1/4} M!$ , asymptotically as  $N$  goes to infinity. The values  $\alpha$  and  $\beta$  are explicit and positive.

*Proof.* The choice of  $u$  in (16) leads to the oscillatory function  $\hat{u}(\xi) = \cos(\frac{\pi}{2}\xi)$ . The modulus  $|\hat{u}|$  is maximal at every point  $\xi \in 2\mathbb{N}$ . Let  $\xi_0 \in 2\mathbb{N}$  and set  $r = \frac{1}{4}$ . For any  $\xi \in [\xi_0 - r, \xi_0 + r]$ , we have

$$\begin{aligned} (|\hat{u}|^2)''(\xi) &= \frac{\pi^2}{2} \left( \sin^2\left(\frac{\pi}{2}\xi\right) - \cos^2\left(\frac{\pi}{2}\xi\right) \right) \\ &\leq -\frac{\pi^2}{2\sqrt{2}}. \end{aligned}$$

Let  $p \in \mathbb{N}$ . The conditions 1 and 2 of Theorem 1 are satisfied with  $Z = 2p\mathbb{N} \cap [-N/2, N/2]$ ,  $K = \lfloor N/(2p) \rfloor$ ,  $r = 1/4$ ,  $c = \frac{\pi^2\sqrt{2}}{8}$ ,  $\delta = 2p$ . Further notice that for every set  $\Xi \in Z^M$ ,  $F(\Xi) = M$ . We can consequently set  $C = M$  for condition 3. For this example, the condition (15) therefore reads

$$M < \frac{\pi^2\sqrt{2}(p - 1/4)}{128} < 0.11(p - 1/4). \quad (17)$$

As long as this condition is satisfied, Theorem 1 allows to conclude on the existence of  $\binom{M}{N/2p} \cdot M!$

maximizers. Taking  $p = \lfloor \sqrt{N} \rfloor$  and  $M = \lfloor 0.11 \cdot \sqrt{N} \rfloor$  yields a number of minimizers larger than  $\binom{\lfloor 0.11 \cdot \sqrt{N} \rfloor}{\lfloor \sqrt{N}/2 \rfloor} \cdot M!$ . Using Stirling formula, we have as  $\gamma$  goes to infinity while  $0 < \theta < 1$

$$\binom{\theta\gamma}{\gamma} \sim \alpha \frac{e^{\beta\gamma}}{\sqrt{\gamma}},$$

with  $\alpha = (2\pi\theta(1 - \theta))^{-1/2}$  and  $\beta = -\theta \log(\theta) - (1 - \theta) \log(1 - \theta)$ . Setting  $\theta = 0.22$  concludes the proof.  $\square$

## 4.2 Numerical illustration of Theorem 1

In this section we illustrate Theorem 1 through numerical examples in Fig. 2. From the left to the right, we used three different 1D signals: a high frequency cosine, a low frequency sine and a Gaussian. We plot the different energy landscapes, for  $M = 2$  measurements at locations  $\Xi = \{\xi_1, \xi_2\}$  and  $N = 16$ . From the top to the bottom, we display  $J_1 = \frac{1}{2}\|u\|_2^2 - F + G_1$ ,  $J_2 = \frac{1}{2}\|u\|_2^2 - F + G_2$ , the functions  $G_1$ ,  $G_2$ ,  $-F$  and the modulus of the Fourier transform  $\xi \mapsto |\hat{u}(\xi)|$ . In order to understand the effect of the signal's structure, the local minima of  $J_1$ ,  $J_2$ ,  $G_1$ ,  $G_2$  and  $-F$  are represented with red dots.

First notice that the cost functions are symmetric with respect to the diagonal. This simply reflects the fact that permutation of points lead to the same energy, and this illustrates the factor  $M!$  in Theorem 1.

As can be seen in all cases, the functions  $G_1$  and  $G_2$  vanish far away from the diagonal (see Corollary 2 below). These point configurations correspond to well-spread sampling schemes. On the contrary, the function  $-F$  can have a large amplitude even outside the diagonal. These two properties are the main ingredients to prove Theorem 1.

The left column (high frequency cosine), corresponds to the example in Corollary 1. We see a number of minimizers that seems quadratic in  $N$  for  $M = 2$ . The center column (low frequency sine) shows that the number of minimizers decreases with a higher regularity of the signal, by reducing the oscillations in  $F$ . On the right (Gaussian function), we illustrate a case where  $F$  has only one local maximum. Even in this case, the function  $J$  still has valleys with shallow local minima. The same phenomenon appears in the center (low frequency sine). Notice that this phenomenon is not captured by Theorem 1, which only relies on local minimizers of  $F$ . In these two examples, the oscillations are induced by the function  $G$ , which we do not explore in this paper.

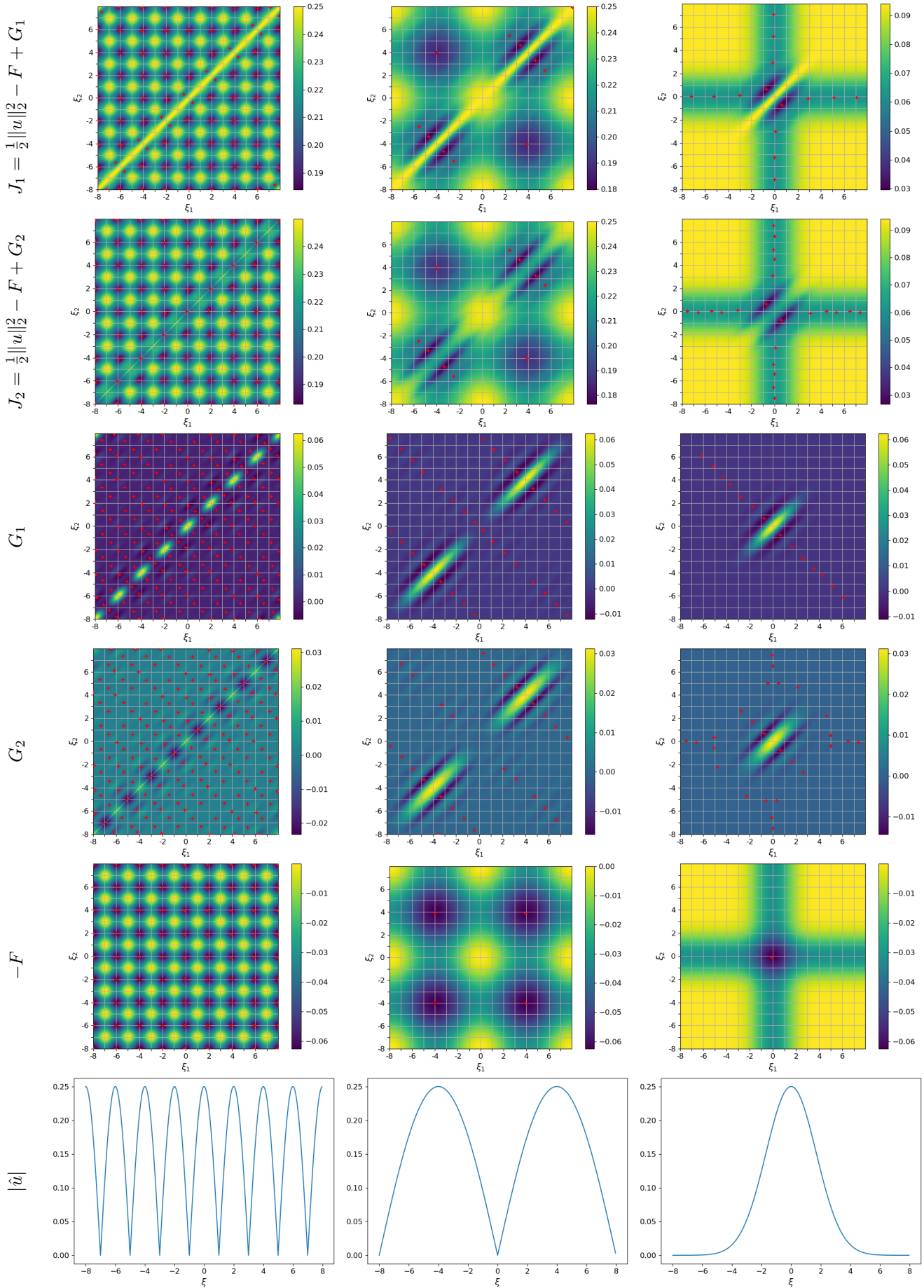


Figure 2: The energy profile for  $M = 2$  and three different signals  $\hat{u}$ : a high frequency cosine, a low frequency sine and a Gaussian (from left to right). From top to bottom, we represent  $J_1$ ,  $J_2$ ,  $G_1$ ,  $G_2$ ,  $F$  and  $|\hat{u}|$ . The red dots represent local minima.

### 4.3 Flatness for high frequencies

In this paragraph we show that the partial derivatives of the cost function may vanish, for indexes corresponding to high frequencies. This explains another practical difficulty in Fourier sampling optimization: without using variable metric techniques, the sampling points located in the high frequencies move very slowly. Though our proof only applies to the function  $J_1$ , this effect also seems to occur for  $J_2$ . See for instance the four corners of Fig. 2, right.

**Proposition 4.** *Letting  $r$  denote the residual error function*

$$r(\Xi) = A(\Xi)A(\Xi)^*u - u, \quad (18)$$

*the gradient of the cost function  $J_1$  reads:*

$$\nabla J_1(\Xi) = \text{Re} \left( \nabla \left( \hat{u}(\Xi) \odot \overline{\hat{r}(\Xi)} \right) \right), \quad (19)$$

*where  $\nabla$  in the right-hand-side denotes the usual derivative in 1D or the gradient in higher dimension and where  $\odot$  is the coordinate-wise (Hadamard) product.*

The proof of Proposition 4 is postponed to Section 7.3.

**Theorem 2** (Vanishing gradients for high frequencies). *Consider a signal  $u \in \mathbb{C}^N$  and a point configuration  $\Xi \in \mathbb{R}^M$ . Under the decay assumptions*

$$|\hat{u}(\xi)| \lesssim \frac{1}{|\xi|^\alpha} \quad \text{and} \quad |\hat{u}'(\xi)| \lesssim \frac{1}{|\xi|^\alpha}, \quad (20)$$

*with  $\alpha > 0$ , we have*

$$\left| \frac{\partial J_1(\Xi)}{\partial \xi_m} \right| \lesssim \frac{\|\hat{u}(\Xi)\|_1}{\text{md}(\Xi)|\xi_m|^\alpha}. \quad (21)$$

The decay assumption appear naturally in the continuous setting, when considering signals  $u$  from Sobolev spaces  $H^k$  with  $k$  derivatives in  $L^2$ .

## 5 Numerical tips and tricks

In this section we propose some solutions to mitigate the issues raised in Section 4 and we illustrate them numerically.

### 5.1 The effect of using a large dataset

In Theorem 1, we proved existence of many local minimizers in the case  $P = 1$ , which corresponds to a unique signal. Let us now assume that we have access to  $P$  signals  $u_1, \dots, u_P$  in  $\mathbb{C}^N$ . The analysis carried out to prove Theorem 1 can

be replicated verbatim. The only difference being that every occurrence of  $|\hat{u}|^2$  must be replaced by  $\rho_P \stackrel{\text{def}}{=} \frac{1}{P} \sum_{p=1}^P |\hat{u}_p|^2$ . The function  $\rho_P$  can be understood as the average power spectral density of the family  $u_1, \dots, u_P$ . As highlighted in Theorem 1, two important factors that can create spurious minimizers are i) the number  $K$  of strict maximizers of  $\rho_P$  and ii) the curvature  $c$  at these maximizers.

As  $P$  increases, we typically expect the density  $\rho_P$  to become smoother. This effect is illustrated for a simple family of shifted and dilated rectangular functions in Fig. 3. As can be seen, both the number of maxima and the curvature  $c$  of  $\rho_P$  in Theorem 1 decay with  $P$ . For  $N = 128$ , we display the average power spectral density for  $P$  ranging from 1 to  $10^3$ . Each signal is defined by

$$u[n] = \int_{n-\frac{1}{2}}^{n+\frac{1}{2}} \mathbb{1}_{[a,b]}(x) dx, \quad (22)$$

where  $a$  and  $b$  are drawn uniformly in the range  $[-N/2 + 1, N/2 - 1]$ . The discrete signals are then renormalized so that  $\|u\|_2 = 1$ .

The same experiment can be reproduced in a more relevant framework from a practical viewpoint. The average power spectral density for 2D knee images of the fastMRI database [14] are represented in Fig. 4. The image are of size  $320 \times 320$ . The local maximizers are computed and displayed with red dots in Fig. 4. In that case, increasing the family size  $P$  reduces the number of maximizers at a slow rate. Indeed they slightly increase from 13k points in the case  $P = 1$  to 14k in the case  $P = 100$  and then start to decrease to 11k for  $P = 10000$ . However, the curvature  $c$  decays much faster. As a conclusion, we see that *using large families of signals can reduce asymptotically the number and the size of the basins of attraction of some spurious minimizers.*

### 5.2 Stochastic gradient descent

When using a large family of signals, the cost function (3) naturally lends itself to the use of stochastic gradient descents (SGD), see [12, 13] that address large MRI datasets. Contrarily to a deterministic gradient descent, which is known to converge to critical points under mild regularity conditions, the stochastic gradient with a fixed step size does not converge. The method is known to end up frolicking in the neighborhood of local critical points [5]. The radius of the neighborhood depends on the stochastic gradient variance and on the step-size. Intuitively, *using stochastic gradients algorithms should therefore allow escaping local minimizers.* We will

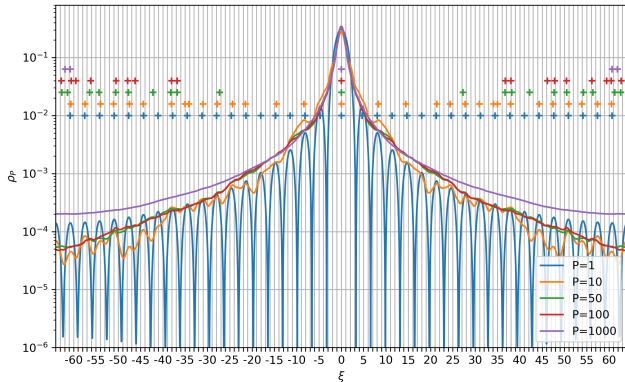


Figure 3: Average power spectral density  $\rho_P$  for families of rectangular functions with different sizes  $P$ . The dots represent local maxima of  $\rho_P$  for different values of  $P$ .

showcase this effect in the forthcoming numerical experiments.

### 5.3 Variable metric

In Section 4.3, Theorem 2 states that the gradient of  $J_1$  might vanish in the high frequency domain. Using second order information is a well known remedy to mitigate this effect. In this work, we propose a simple method which corresponds to a variable diagonal metric with well-chosen coefficients.

As shown in Theorem 2, the gradient vanishes with a rate depending on the Fourier transform magnitude  $|\hat{u}|$ . For a dataset, this decay is somewhat captured by the average power spectral density  $\rho_P(\xi) \stackrel{\text{def}}{=} \frac{1}{P} \sum_{p=1}^P |\hat{u}_p(\xi)|^2$ . Hence, we propose to compute  $\rho_P$  once and for all on a fine grid ( $20 \times N$  discretization points in our example). The function  $\rho_P$  is then linearly interpolated in between the grid points during the gradient descent. At each gradient iteration we replace  $\frac{\partial J_1(\Xi)}{\partial \xi_m}$  by

$$\frac{1}{\rho_P(\xi_m)^\beta} \frac{\partial J_1(\Xi)}{\partial \xi_m}, \quad (23)$$

where  $\beta$  is a constant that has to be set empirically. From numerical experiments  $\beta \in [1, 2]$  shows good performance. In all the experiments presented hereafter we use  $\beta = 1$ . We will see later in the numerical experiments, that *this variable metric significantly accelerates the convergence for sampling points located in high frequencies*.

### 5.4 Numerical illustrations

In this section, we aim at illustrating numerically the different results established previously. We aim at reconstructing 1D signals of size  $N = 128$  from

$M = 64$  measurements in the Fourier domain. We suppose that  $P$  rectangular signals generated using (22) are given. We illustrate our findings with the back-projection reconstructor associated to the cost function  $J_1$ , but similar results have been obtained with the pseudo-inverse. As we are working in 1D with small dimensions  $N$  and  $M$ , at each iteration, the whole matrix  $A(\Xi)^*$  is evaluated and the gradient  $\nabla J_1$  is computed directly from the analytic expression (19). We first use a fixed step gradient descent algorithm in order to showcase the convergence dynamics of the algorithm. The initialization of  $\Xi$  is a subgrid with a constant spacing of 2. The following experiments are conducted:

**Effect of the dataset size  $P$**  We first vary the number of signals by taking  $P = 1$  and  $P = 1000$ . The evolution of  $\Xi$  is displayed in Fig. 5, respectively top-left and top-center. The history of the cost function is given in Fig. 6. For this experiment, we expect that a good sampling scheme consists of low frequencies sampled at the Shannon-Nyquist rate. In this regard, the sampling scheme obtained in Fig. 5 for  $P = 1000$  is more satisfactory than the one obtained for  $P = 1$ . In the case  $P = 1000$ , the displacement of  $\Xi$  is more important, suggesting that some local minima have been discarded.

**Variable metric** We then study, for  $P = 1000$  the effect of a variable metric gradient descent as described in Section 5.3. We also compare this approach to an L-BFGS algorithm with a line search and with a Hessian estimated using the last 8 gradients. In Fig. 5, the usual gradient algorithm is at the top-center, the variable metric gradient descent is at the bottom-center and the L-BFGS algorithm is at the bottom-left. The cost function evolution is displayed in Fig. 6. Using a variable metric results in a huge speed-up of the algorithm. This is particularly visible for points  $\xi$  located at high frequencies, which is another illustration of Theorem 2. For this example, the L-BFGS algorithm converges slightly faster than the variable metric gradient descent in the early iterations. However, its per-iteration cost is much higher since it uses a line search and a non diagonal metric. Since the L-BFGS algorithm can be seen as a state-of-the-art quasi-Newton method, the proposed empirical metric (23) seems remarkably efficient.

**Stochastic gradient descent** Finally in Fig. 5 right column, we investigate the use of a fixed-step stochastic gradient descent algorithm with a batch size of 1. In that experiment, a new ran-



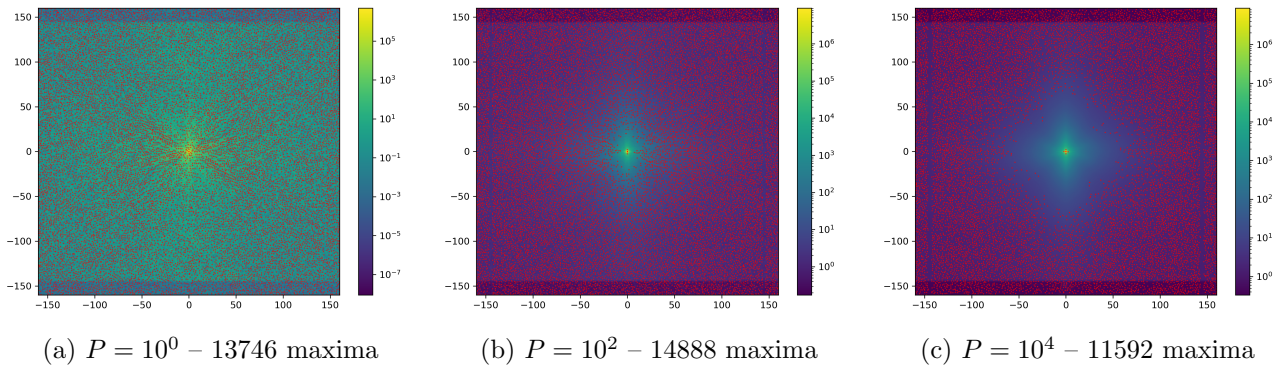


Figure 4: Average power spectral density  $\rho_P$  for a subset of images from the knee dataset of fastMRI. The image size is  $N = 320$  and the red dots represent local maximizers.

dom signal is generated at every iteration using the model (22) and the stochastic gradient is computed with respect to that signal only. The trajectory of the vanilla SGD is comparable with the one obtained using a deterministic gradient descent for  $P = 1000$  in Fig. 5 top-center. The variable metric trick significantly improves the convergence speed and more importantly, the final points configuration. As a conclusion, the variable metric SGD algorithm seems to be able to escape spurious minimizers and to take advantage of the averaging effect of the large dataset without the struggle of computing the gradient over a large dataset.

**Comparison of the sampling schemes** The final sampling schemes are not directly comparable in terms of cost function because the objective function is computed over different datasets. In Table 1, we therefore report the cost function computed on a specific set of signals. This set contains the  $P = 1000$  signals that are used in the numerical illustrations of Fig. 5 center column. When tested against a large dataset, the final configuration obtained for  $P = 1$  seems highly sub-optimal. This effect is most likely due to a convergence to a local minimizer and also to the fact that the sampling scheme is not adapted to a whole family but only to a single signal. The remarkable observation that can be made from Table 1 is that the optimal configuration obtained with the variable metric SGD performs better on the dataset of  $P = 1000$  signals than the experiment conducted in Fig. 5 which is tailored for this dataset. This shows that the usual deterministic algorithms are stuck in local minima even with large datasets. On the contrary, the variable metric SGD algorithm seems effective.

These numerical results highlight the effectiveness of the different tricks suggested in this section: the use of a variable metric to handle high frequencies and a stochastic optimization to avoid

local minima.

## 6 Conclusion

We highlighted two obstacles to the convergence of gradient based algorithms for sampling schemes optimization. The first one is a high number of local minimizers and the second one is a vanishing gradient phenomenon for high frequencies. The first obstruction can be mitigated with a regularization by averaging and the second one by an adhoc variable metric gradient descent. Unfortunately, these tricks still seem insufficient to avoid some local minimizers. The additional use of a stochastic gradient instead of a deterministic gradient approach seems to leverage most difficulties in a simplified 1D setting.

Many state-of-the-art reconstructors are based on a quadratic data fidelity term and we expect that some of the techniques used in this paper in the linear case can be reused even in a nonlinear setting. This is left for future research.

## 7 Proofs

### 7.1 Proof of Proposition 3

*Proof.* Let us start with  $J_1$ . We develop the squared norm, leading to :

$$J_1(\Xi) = \frac{1}{2} \|u\|_2^2 + \frac{1}{2} \|A(\Xi)A(\Xi)^*u\|_2^2 - \text{Re}(\langle A(\Xi)A(\Xi)^*u, u \rangle).$$

Now, observe that

$$\begin{aligned} \|A(\Xi)A(\Xi)^*u\|_2^2 &= \|A(\Xi)\hat{u}(\Xi)\|_2^2 \\ &= \langle A(\Xi)\hat{u}(\Xi), A(\Xi)\hat{u}(\Xi) \rangle = \langle A(\Xi)^*A(\Xi)\hat{u}(\Xi), \hat{u}(\Xi) \rangle. \end{aligned}$$

Then, (13) is a direct consequence of

$$\langle A(\Xi)A(\Xi)^*u, u \rangle = \langle A(\Xi)^*u, A(\Xi)^*u \rangle = \|\hat{u}(\Xi)\|_2^2.$$

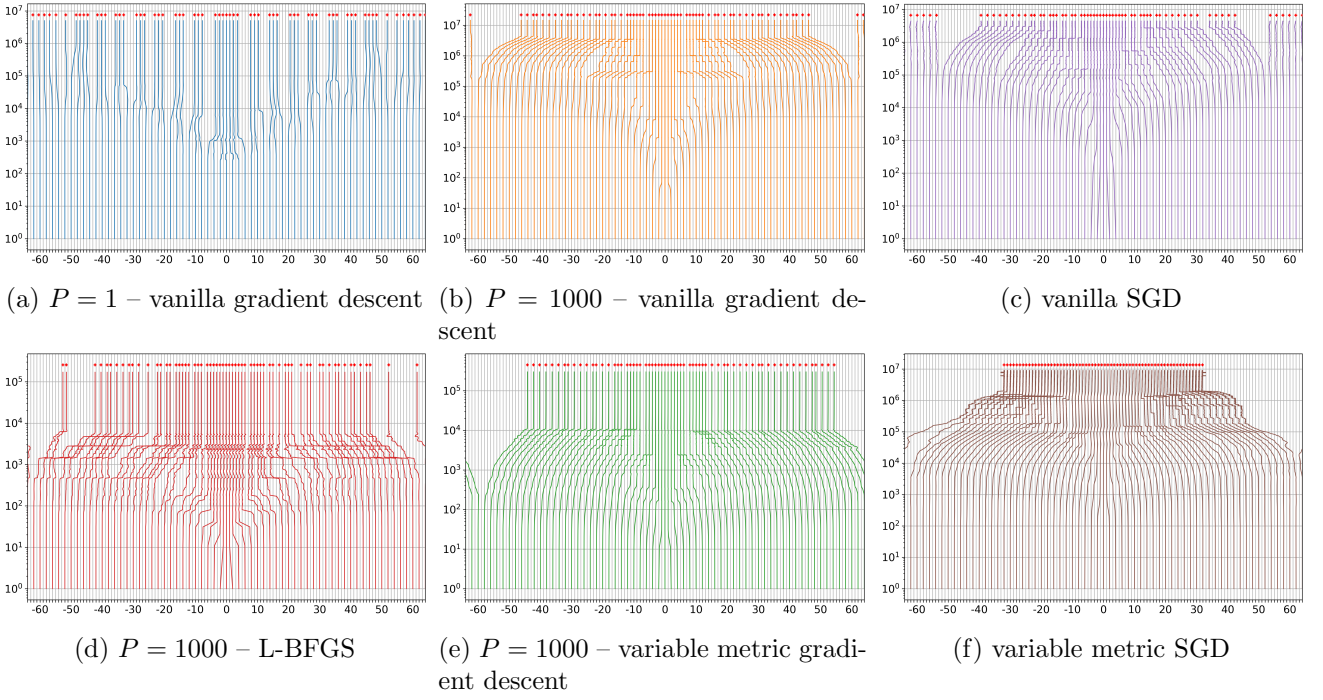


Figure 5: Trajectories of  $\Xi$  the back-projection reconstructor  $J_1$  and a fixed-step gradient descent. The iterations are represented on the vertical axis, and the horizontal axis corresponds to  $\xi$  and is periodic. The initialization is a uniform subgrid and is seen on the axis  $y = 0$  of the top and middle figures. Left and center: trajectories of  $\Xi$  for different sizes of signals families. The objective function is given in Fig. 6. The right column represents trajectories of  $\Xi$  using a stochastic gradient descent with one signal in the batch that is different at each iteration. The trajectories in the stochastic case have been averaged over the last 10000 iterations.

Test case	$P = 1$	$P = 1000$	$P = 1000$ with var. metric	L-BFGS	SGD	SGD with var. metric
Eff.	$9.07 \times 10^{-2}$	$2.68 \times 10^{-2}$	$2.38 \times 10^{-2}$	$2.41 \times 10^{-2}$	$6.63 \times 10^{-2}$	$1.00 \times 10^{-2}$

Table 1: Effectiveness of the sampling schemes obtained with different strategies on a dataset of 1000 signals. The table contains the average reconstruction error  $J_1$  over the dataset. This dataset is the one used in the case  $P = 1000$ , see Fig. 5 center column.

Now, let us turn to  $J_2$ . Using Pythagorean theorem and the fact that  $(A(\Xi)^*)^+ A(\Xi)^* = \Pi_{\text{ran}(A(\Xi))}$  we have:

$$\begin{aligned}
J_2(\Xi) &= \frac{1}{2} \|(A(\Xi)^*)^+ A(\Xi)^* u - u\|_2^2 \\
&= \frac{1}{2} \|\Pi_{\text{ran}(A(\Xi))}(u) - u\|_2^2 \\
&= \frac{1}{2} \|u\|_2^2 - \frac{1}{2} \|\Pi_{\text{ran}(A(\Xi))}(u)\|_2^2
\end{aligned}$$

Then, using the identity

$$\Pi_{\text{ran}(A(\Xi))} = A(\Xi)(A(\Xi)^* A(\Xi))^+ A(\Xi)^*,$$

we obtain

$$\begin{aligned}
\frac{1}{2} \|\Pi_{\text{ran}(A(\Xi))}(u)\|_2^2 &= \frac{1}{2} \langle u, \Pi_{\text{ran}(A(\Xi))}(u) \rangle \\
&= \frac{1}{2} \langle A(\Xi)^* u, (A(\Xi)^* A(\Xi))^+ A(\Xi)^* u \rangle \\
&= \frac{1}{2} \langle \hat{u}(\Xi), (A(\Xi)^* A(\Xi))^+ \hat{u}(\Xi) \rangle.
\end{aligned}$$

Adding and subtracting  $\frac{1}{2} \|\hat{u}(\Xi)\|_2^2$  finishes the proof of (14).  $\square$

## 7.2 Proof of Theorem 1

### 7.2.1 Controlling the amplitude of $G$

Significant progress have been made lately in the control of the extreme eigenvalues of Vandermonde matrices, which play a pivotal role in algebraic number theory [4, 9, 3, 1]. The tightest results for well separated schemes was recently obtained in [1]. Rewriting their result in our formalism, we obtain the following inequality.

**Proposition 5** (Conditioning of Vandermonde matrices [1]). *Let  $\Xi = (\xi_1, \dots, \xi_M)$  denote a set of distinct sampling points. Let  $L(\Xi) = A(\Xi)A(\Xi)^*$ .*

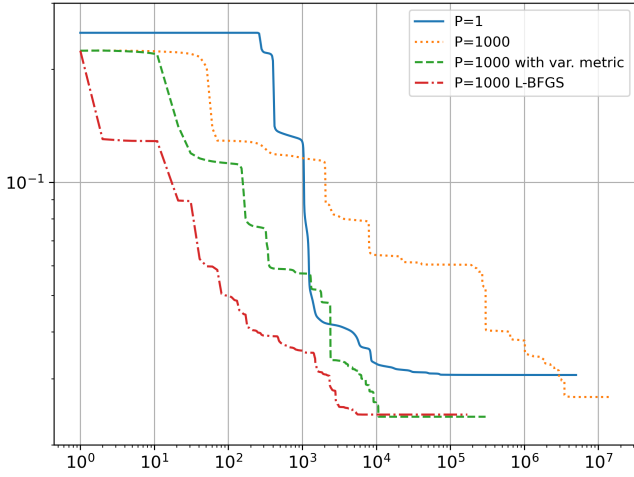


Figure 6: Objective function  $J_1$  (back-projection) for the different experiments in Fig. 5 in the deterministic case.

The following inequality holds

$$\left(1 - \frac{1}{\text{md}(\Xi)}\right) \text{Id} \preceq L(\Xi) \preceq \left(1 + \frac{1}{\text{md}(\Xi)}\right) \text{Id}. \quad (24)$$

*Proof.* This is a direct consequence of [1, eq. (31)] up to renormalizations.  $\square$

The above result allows to obtain the following corollary.

**Corollary 2** (Function  $G$  is small for well separated schemes). *When all the sampling points  $\xi_m$  are distinct, we have*

$$|G_1(\Xi)| \leq \frac{F(\Xi)}{\text{md}(\Xi)} \quad (25)$$

and for  $\text{md}(\Xi) > 1$ ,

$$|G_2(\Xi)| \leq \frac{F(\Xi)}{\text{md}(\Xi) - 1} \quad (26)$$

*Proof.* We have

$$\begin{aligned} |G_1(\Xi)| &= \left| \frac{1}{2} \langle (A(\Xi)^* A(\Xi) - \text{Id}) \hat{u}(\Xi), \hat{u}(\Xi) \rangle \right| \\ &\leq \frac{1}{2} \| (A(\Xi)^* A(\Xi) - \text{Id}) \|_{2 \rightarrow 2} \|\hat{u}(\Xi)\|_2^2 \\ &\leq \| (A(\Xi)^* A(\Xi) - \text{Id}) \|_{2 \rightarrow 2} F(\Xi) \leq \frac{F(\Xi)}{\text{md}(\Xi)}. \end{aligned}$$

For  $\Xi$  s.t.  $\text{md}(\Xi) > 1$ ,  $L(\Xi)$  is invertible so that  $L(\Xi)^+ = L(\Xi)^{-1}$ . Proposition 5 then yields

$$\frac{\text{md}(\Xi)}{1 + \text{md}(\Xi)} \text{Id} \preceq L(\Xi)^{-1} \preceq \frac{\text{md}(\Xi)}{\text{md}(\Xi) - 1} \text{Id}, \quad (27)$$

which implies

$$-\frac{1}{1 + \text{md}(\Xi)} \text{Id} \preceq L(\Xi)^{-1} - \text{Id} \preceq \frac{1}{\text{md}(\Xi) - 1} \text{Id},$$

and  $\|\text{Id} - L(\Xi)^{-1}\|_{2 \rightarrow 2} \leq \frac{1}{\text{md}(\Xi) - 1}$ . It suffices to apply the same reasoning as for  $G_1$  to conclude.  $\square$

This inequality tells us that for well separated sampling schemes, the functions  $G_1$  and  $G_2$  cannot be too large. For instance, consider  $M$  roughly equispaced points on  $[-N/2, N/2]$ . Then  $\text{md}(\Xi) \simeq N/M$  and  $|G_1(\Xi)| = O\left(\frac{M}{N} \|\hat{u}(\Xi)\|_2^2\right)$ .

## 7.2.2 Combining the previous results

In what follows,  $G$  represents either  $G_1$  or  $G_2$ . Under the hypotheses of Theorem 1, first notice that any set  $\Xi \in Z^M$  is a local maximizer of  $F$ . Indeed any perturbation of the individual sampling locations  $\xi_m$  results in a decay of the captured energy.

There are  $\binom{M}{K}$  possible sampling configurations when all the points belong to  $Z$ . Let  $\bar{\Xi} = \{\bar{\xi}_1, \dots, \bar{\xi}_M\}$  denote one of them. The idea of the proof is to show that there is a local maximizer of  $F - G$  in the following neighborhood  $B = [\bar{\xi}_1 - r, \bar{\xi}_1 + r] \times \dots \times [\bar{\xi}_M - r, \bar{\xi}_M + r]$ . A sufficient condition for the set  $B$  to contain a local maximizer of  $F - G$  is that  $F(\bar{\Xi}) - G(\bar{\Xi}) > F(\Xi) - G(\Xi)$  for all  $\Xi \in \partial B$  (the boundary of  $B$ ) since  $F - G$  is continuous.

Using Corollary 2, first notice that

$$\begin{aligned} F(\bar{\Xi}) - G_1(\bar{\Xi}) &\geq F(\bar{\Xi}) - \frac{1}{\text{md}(\bar{\Xi})} F(\bar{\Xi}) \\ &\geq F(\bar{\Xi}) \left(1 - \frac{1}{\text{md}(\bar{\Xi})}\right) \geq F(\bar{\Xi}) \left(1 - \frac{1}{\delta}\right). \end{aligned}$$

For all  $\Xi \in B$ , we have by strict concavity of  $|\hat{u}|$  around each  $\bar{\xi}_m$ ,  $F(\Xi) \leq F(\bar{\Xi}) - \frac{c}{2} \sum_{m=1}^M (\xi_m - \bar{\xi}_m)^2$ . Hence, for sampling sets  $\Xi \in \partial B$  on the boundary of  $B$ , we have

$$F(\Xi) \leq F(\bar{\Xi}) - \frac{cr^2}{2}. \quad (28)$$

In addition for  $\Xi \in \partial B$ ,  $\text{md}(\Xi) \geq \delta - 2r$  and using Corollary 2 again, we obtain:

$$\begin{aligned} F(\Xi) - G_1(\Xi) &\leq F(\bar{\Xi}) - \frac{cr^2}{2} + F(\bar{\Xi}) \frac{1}{\delta - 2r} \\ &= F(\bar{\Xi}) \left(1 + \frac{1}{\delta - 2r}\right) - \frac{cr^2}{2}. \end{aligned}$$

Therefore, the condition

$$F(\bar{\Xi}) \left(1 + \frac{1}{\delta - 2r}\right) - \frac{cr^2}{2} < F(\bar{\Xi}) \left(1 - \frac{1}{\delta}\right) \quad (29)$$

suffices to conclude on the existence of a maximizer of  $J_1$  in the interior of  $B$ . This condition is satisfied for  $F(\bar{\Xi}) < \frac{cr^2(\delta - 2r)\delta}{4(\delta - r)}$  and a fortiori for

$$F(\bar{\Xi}) < \frac{cr^2(\delta - 2r)}{4} \quad (30)$$

The multiplicative factor  $M!$  is related to the fact that for a given maximizer, all the possible permutations of indices give rise to different maximizers.

The same reasoning applies verbatim to  $J_2$  by replacing  $\delta$  with  $\delta - 1$ .

### 7.3 Proof of Proposition 4

*Proof.* Let us consider a point configuration  $\Xi \in \mathbb{R}^M$  and a perturbation  $\epsilon \in \mathbb{R}^M$ . Given a vector of measurements  $\hat{u}(\Xi) \in \mathbb{C}^M$ , we let  $\nabla \hat{u}(\Xi) =$

$\begin{pmatrix} \hat{u}'(\xi_1) \\ \vdots \\ \hat{u}'(\xi_M) \end{pmatrix}$  denote the vector of derivatives at the

sampling locations. Elementary calculus leads to the following identities for every  $\epsilon$  direction of variation:

$$\begin{aligned} (\text{Jac}_A(\Xi)\epsilon)^* &= \text{Jac}_{A^*}(\Xi)\epsilon \\ \nabla \hat{u}(\Xi) \odot \epsilon &= \text{Jac}_{A^*}(\Xi)\epsilon \cdot u. \end{aligned}$$

Then, we apply standard calculus of variations:

$$\begin{aligned} J_1(\Xi + \epsilon) &= J_1(\Xi) + \text{Re}\langle \text{Jac}_A(\Xi)\epsilon \cdot \hat{u}(\Xi), r(\Xi) \rangle \\ &\quad + \text{Re}\langle A(\Xi)\text{Jac}_{A^*}(\Xi)\epsilon \cdot u, r(\Xi) \rangle + o(\|\epsilon\|_2^2) \\ &= J_1(\Xi) + \text{Re}\langle \hat{u}(\Xi), (\text{Jac}_A(\Xi)\epsilon)^* r(\Xi) \rangle \\ &\quad + \text{Re}\langle \nabla \hat{u}(\Xi) \odot \epsilon, \hat{r}(\Xi) \rangle + o(\|\epsilon\|_2^2) \\ &= J_1(\Xi) + \text{Re}\langle \hat{u}(\Xi), \nabla \hat{r}(\Xi) \odot \epsilon \rangle \\ &\quad + \text{Re}\langle \epsilon, \overline{\nabla \hat{u}(\Xi)} \odot \hat{r}(\Xi) \rangle + o(\|\epsilon\|_2^2) \\ &= J_1(\Xi) + \text{Re}\langle \overline{\nabla \hat{r}(\Xi)} \odot \hat{u}(\Xi), \epsilon \rangle \\ &\quad + \text{Re}\langle \epsilon, \overline{\nabla \hat{u}(\Xi)} \odot \hat{r}(\Xi) \rangle + o(\|\epsilon\|_2^2). \end{aligned}$$

Hence, by identification

$$\begin{aligned} \nabla J_1(\Xi) &= \text{Re}\left(\overline{\nabla \hat{r}(\Xi)} \odot \hat{u}(\Xi) + \nabla \hat{u}(\Xi) \odot \overline{\hat{r}(\Xi)}\right) \\ &= \text{Re}\left(\nabla\left(\hat{u}(\Xi) \odot \overline{\hat{r}(\Xi)}\right)\right). \end{aligned}$$

### 7.4 Proof of Theorem 2

By Proposition 4, we have

$$\left| \frac{\partial J_1(\Xi)}{\partial \xi_m} \right| \leq |\hat{u}'(\xi_m)| \cdot |\hat{r}(\xi_m)| + |\hat{u}(\xi_m)| \cdot |\hat{r}'(\xi_m)|.$$

By definition, we have  $\hat{r}(\Xi) = (L(\Xi) - \text{Id})\hat{u}(\Xi)$ , hence

$$|\hat{r}(\xi_m)| \leq \|\hat{r}(\Xi)\|_2 \leq \frac{\|\hat{u}(\Xi)\|_2}{\text{md}(\Xi)}, \quad (31)$$

where we used Proposition 5 to obtain the last inequality. Now, we also wish to control  $|\hat{r}'(\xi_m)|$ . To

this end, first notice that

$$\begin{aligned} \hat{r}'(\xi_m) &= \sum_{m'=1}^M \left( \frac{\partial L(\Xi)_{m,m'}}{\partial \xi_m} \hat{u}(\xi_{m'}) \right. \\ &\quad \left. + L(\Xi)_{m,m'} \hat{u}'(\xi_{m'}) \mathbb{1}_{m=m'} \right) - \hat{u}'(\xi_m) \\ &= \sum_{m'=1}^M \frac{\partial L(\Xi)_{m,m'}}{\partial \xi_m} \hat{u}(\xi_{m'}). \end{aligned}$$

We start with an analytical expression of the matrix  $L(\Xi)$ .

**Proposition 6** (The expression of  $A^*A$ ). *Let  $L(\Xi) \stackrel{\text{def}}{=} A(\Xi)^*A(\Xi)$ . We have*

$$[L(\Xi)]_{m,m'} = \begin{cases} 1 & \text{if } m = m', \\ \frac{1}{N} \exp\left(\frac{\iota\pi(\xi_m - \xi_{m'})}{N}\right) & \\ \quad \times \frac{\sin(\pi(\xi_m - \xi_{m'}))}{\sin\left(\frac{\pi(\xi_m - \xi_{m'})}{N}\right)} & \text{otherwise.} \end{cases} \quad (32)$$

*Proof.* We have:

$$\begin{aligned} [L(\Xi)]_{m,m'} &= \frac{1}{N} \sum_n e^{2\iota\frac{\pi}{N}\langle \xi_{m'} - \xi_m, n \rangle} \\ &= \frac{1}{N} e^{-\iota\pi(\xi_{m'} - \xi_m)} \frac{1 - e^{2\iota\pi(\xi_{m'} - \xi_m)}}{1 - e^{2\iota\frac{\pi}{N}(\xi_{m'} - \xi_m)}} \\ &= \frac{1}{N} e^{-\iota\pi(\xi_{m'} - \xi_m)} \frac{e^{\iota\pi(\xi_{m'} - \xi_m)}}{e^{\iota\frac{\pi}{N}(\xi_{m'} - \xi_m)}} \\ &\quad \times \frac{e^{-\iota\pi(\xi_{m'} - \xi_m)} - e^{\iota\pi(\xi_{m'} - \xi_m)}}{e^{-\iota\frac{\pi}{N}(\xi_{m'} - \xi_m)} - e^{\iota\frac{\pi}{N}(\xi_{m'} - \xi_m)}} \\ &= \frac{1}{N} e^{-\iota\frac{\pi}{N}(\xi_{m'} - \xi_m)} \frac{\sin(\pi(\xi_{m'} - \xi_m))}{\sin\left(\frac{\pi}{N}(\xi_{m'} - \xi_m)\right)}. \end{aligned}$$

□

Now, we will use the following lemma.

**Lemma 1.** *The following bound holds:*

$$\left| \frac{\partial L(\Xi)_{m,m'}}{\partial \xi_m} \right| \leq \frac{\pi}{N} + \frac{4}{\text{dist}(\xi_{m'}, \xi_m)} \leq \frac{\pi}{N} + \frac{4}{\text{md}(\Xi)}.$$

□ *Proof.* Letting  $\delta = \xi_m - \xi_{m'}$ , we have

$$\begin{aligned} \frac{\partial L(\Xi)_{m,m'}}{\partial \xi_m} &= \frac{\pi}{N^2} \times \frac{\iota e^{\iota\frac{\pi}{N}\delta} \sin(\pi\delta)}{\sin\left(\frac{\pi}{N}\delta\right)} \\ &\quad + \frac{\pi}{N} \times \frac{e^{-\iota\frac{\pi}{N}\delta}}{\sin\left(\frac{\pi}{N}\delta\right)} \left( \cos(\pi\delta) - \frac{\sin(\pi\delta)}{N} \times \frac{\cos\left(\frac{\pi}{N}\delta\right)}{\sin\left(\frac{\pi}{N}\delta\right)} \right). \end{aligned}$$

Without loss of generality we consider the case  $0 \leq \delta \leq N/2$ . Using  $\left| \frac{\sin(\pi\delta)}{N \sin\left(\frac{\pi}{N}\delta\right)} \right| \leq 1$  let us remark that

$$\begin{aligned} \left| \frac{\partial L(\Xi)_{m,m'}}{\partial \xi_m} \right| &\leq \frac{\pi}{N} \\ &\quad + \frac{\pi}{N} \left| \frac{1}{\sin\left(\frac{\pi}{N}\delta\right)} \left( \frac{\sin(\pi\delta) \cos\left(\frac{\pi}{N}\delta\right)}{N \sin\left(\frac{\pi}{N}\delta\right)} - \cos(\pi\delta) \right) \right|. \end{aligned}$$

Using the inequality  $\left| \frac{\sin(\pi\delta)}{N \sin(\frac{\pi}{N}\delta)} \right| \leq 1$  again, we obtain

$$\left| \frac{\sin(\pi\delta) \cos(\frac{\pi}{N}\delta)}{N \sin(\frac{\pi}{N}\delta)} - \cos(\pi\delta) \right| \leq \left| \cos(\frac{\pi}{N}\delta) \right| + 1 \leq 2.$$

Finally, using the inequality  $\sin(x) \geq x/2$  for  $x \in (0, \pi/2)$ , we get  $\left| \frac{\partial L(\Xi)_{m',m}}{\partial \xi_{m'}} \right| \leq \frac{\pi}{N} + \frac{4}{\delta}$ .  $\square$

Lemma 1 and a Cauchy-Schwarz inequality provides the following bound:

$$|\hat{r}'(\xi_m)| \leq \left( \frac{\pi}{N} + \frac{4}{\text{md}(\Xi)} \right) \|\hat{u}(\Xi)\|_1.$$

Combining everything finally yields:

$$\left| \frac{\partial J_1(\Xi)}{\partial \xi_m} \right| \leq |\hat{u}'(\xi_m)| \cdot \frac{\|\hat{u}(\Xi)\|_2}{\text{md}(\Xi)} + |\hat{u}(\xi_m)| \cdot \|\hat{u}(\Xi)\|_1 \cdot \left( \frac{\pi}{N} + \frac{4}{\text{md}(\Xi)} \right).$$

Under the decay assumptions of Theorem 2, we obtain

$$\left| \frac{\partial J_1(\Xi)}{\partial \xi_m} \right| \lesssim \frac{\|\hat{u}(\Xi)\|_1}{\text{md}(\Xi) |\xi_m|^\alpha}.$$

## References

- [1] Céline Aubel and Helmut Bölcskei. Vandermonde matrices with nodes in the unit disk and the large sieve. *Applied and Computational Harmonic Analysis*, 47(1):53–86, 2019.
- [2] Cagla Deniz Bahadir, Adrian V Dalca, and Mert R Sabuncu. Learning-based optimization of the under-sampling pattern in mri. In *International Conference on Information Processing in Medical Imaging*, pages 780–792. Springer, 2019.
- [3] Dmitry Batenkov, Laurent Demanet, Gil Goldman, and Yosef Yomdin. Conditioning of partial nonuniform fourier matrices with clustered nodes. *SIAM Journal on Matrix Analysis and Applications*, 41(1):199–220, 2020.
- [4] Enrico Bombieri. On the large sieve. In *Goldbach Conjecture*, pages 227–252. World Scientific, 1984.
- [5] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [6] Alban Gossard, Frédéric de Gournay, and Pierre Weiss. Off-the-grid data-driven optimization of sampling schemes in mri. *arXiv preprint arXiv:2010.01817*, 2020.
- [7] Baran Gözcü, Rabeeh Karimi Mahabadi, Yen-Huan Li, Efe Ilıcak, Tolga Cukur, Jonathan Scarlett, and Volkan Cevher. Learning-based compressive mri. *IEEE transactions on medical imaging*, 37(6):1394–1406, 2018.
- [8] Kyong Hwan Jin, Michael Unser, and Kwang Moo Yi. Self-supervised deep active accelerated mri. *arXiv preprint arXiv:1901.04547*, 2019.
- [9] Ankur Moitra. Super-resolution, extremal functions and the condition number of vandermonde matrices. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 821–830, 2015.
- [10] Alan Oppenheim, Don Johnson, and Kenneth Steiglitz. Computation of spectra with unequal resolution using the fast fourier transform. *Proceedings of the IEEE*, 59(2):299–301, 1971.
- [11] Ferdia Sherry, Martin Benning, Juan Carlos De los Reyes, Martin J Graves, Georg Maierhofer, Guy Williams, Carola-Bibiane Schönlieb, and Matthias J Ehrhardt. Learning the sampling pattern for mri. *IEEE Transactions on Medical Imaging*, 39(12):4310–4321, 2020.
- [12] Guanhua Wang, Tianrui Luo, Jon-Fredrik Nielsen, Douglas C Noll, and Jeffrey A Fessler. B-spline parameterized joint optimization of reconstruction and k-space trajectories (bjork) for accelerated 2d mri. *IEEE Transactions on Medical Imaging*, 2021.
- [13] Tomer Weiss, Ortal Senouf, Sanketh Vedula, Oleg Michailovich, Michael Zibulevsky, and Alex Bronstein. Pilot: Physics-informed learned optimal trajectories for accelerated mri. *arXiv preprint arXiv:1909.05773*, 2019.
- [14] Jure Zbontar, Florian Knoll, Anuroop Sriram, Tullie Murrell, Zhengnan Huang, Matthew J Muckley, Aaron Defazio, Ruben Stern, Patricia Johnson, Mary Bruno, et al. fastmri: An open dataset and benchmarks for accelerated mri. *arXiv preprint arXiv:1811.08839*, 2018.
- [15] Marcelo VW Zibetti, Gabor T Herman, and Ravinder R Regatte. Fast data-driven learning of mri sampling pattern for large scale problems. *arXiv preprint arXiv:2011.02322*, 2020.