The troublesome kernel: why deep learning for inverse problems is typically unstable

Nina M. Gottschling* Vegard Antun[†] Ben Adcock[‡] Anders C. Hansen[§]

January 7, 2020

Abstract

There is overwhelming empirical evidence that Deep Learning (DL) leads to unstable methods in applications ranging from image classification and computer vision to voice recognition and automated diagnosis in medicine. Recently, a similar instability phenomenon has been discovered when DL is used to solve certain problems in computational science, namely, inverse problems in imaging. In this paper we present a comprehensive mathematical analysis explaining the many facets of the instability phenomenon in DL for inverse problems. Our main results not only explain why this phenomenon occurs, they also shed light as to why finding a cure for instabilities is so difficult in practice. Additionally, these theorems show that instabilities are typically not rare events – rather, they can occur even when the measurements are subject to completely random noise – and consequently how easy it can be to destablise certain trained neural networks. We also examine the delicate balance between reconstruction performance and stability, and in particular, how DL methods may outperform state-of-the-art sparse regularization methods, but at the cost of instability. Finally, we demonstrate a counterintuitive phenomenon: training a neural network may generically not yield an optimal reconstruction method for an inverse problem.

Keywords: Deep learning, stability, inverse problems, imaging, sparse regularization

Mathematics Subject Classification (2010): 65R32, 94A08, 68T05, 65M12

1 Introduction

It is impossible to overstate the impact that Deep Learning (DL) has had in recent years in core machine learning applications such as image classification, speech recognition and natural language processing. Perhaps unsurprisingly, the development and use of DL for challenging problems in the computational sciences has recently become an active area of inquiry. Areas of particular note include numerical PDEs [17, 48], discovering PDE dynamics [40], Uncertainty Quantification and high-dimensional approximation [42].

Arguably, however, the area of computational science in which DL has been most actively investigated over the last several years is inverse problems in imaging. Image reconstruction from measurements is an important task in a wide range of scientific, industrial and medical applications, including, but by no means limited to, electron and fluorescence microscopy, seismic imaging, Nuclear Magnetic Resonance (NMR), Magnetic Resonance Imaging (MRI) and X-Ray CT. The last several years have witnessed the emergence of a variety of different trained Neural Networks (NNs) for image reconstruction which claim to achieve competitive, and sometimes

^{*}University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK (nmg43@cam.ac.uk)

[†]University of Oslo, P.O box 1053, Blindern, 0316 Oslo, Norway (vegarant@math.uio.no)

[‡]Simon Fraser University, 8888 University Drive Burnaby, BC V5A 1S6, Canada (ben_adcock@sfu.ca)

[§]University of Cambridge, Wilberforce Road, Cambridge CB3 0WA, UK (ach70@cam.ac.uk)

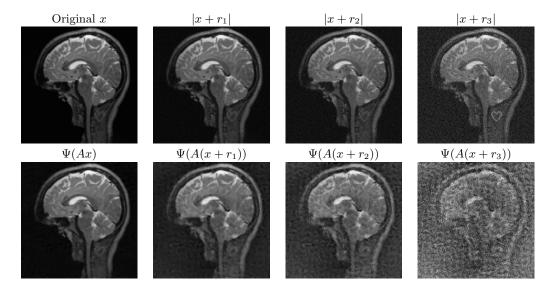


Figure 1: **(Unstable neural network).** The effect of small perturbations on the AUTOMAP network [50] for recovering an image x from its measurements y = Ax. Here $A \in \mathbb{C}^{m \times N}$ is a subsampled discrete Fourier transform, which is the standard mathematical model for MRI. The left column shows that the network, denoted by Ψ , recovers the original image x well. However, as shown in the second to fourth columns, small perturbations $||r_1|| < ||r_2|| < ||r_3||$ of x cause large artefacts in the recovered images $\Psi(A(x+r_i))$. This experiment is from [6].

even superior, performance to current state-of-the-art techniques [7]. Notably, their potential has been described by *Nature* as 'transformative' [44].¹

1.1 Accuracy and stability in computational science

As trained NN algorithms begin to percolate into the computational scientist's toolkit, it is important that they be examined through the traditional pillars of numerical analysis, namely, accuracy and stability. There is a growing awareness that such techniques have not yet been subject to the same rigorous standards as more well-established methods in scientific computing [8]. Moreover, there is evidence that such techniques, in their current guise at least, do not yet meet these standards. For instance in image reconstruction, it has recently been demonstrated that existing DL algorithms, despite offering purportedly 'superior immunity to noise' [44], are often highly unstable [6] (see also [27]). To be precise, a small perturbation in the true image or its measurements may cause severe artefacts in the reconstruction. Fig. 1 gives an illustration of this so-called instability phenomenon. By contrast, as is also shown in this figure, existing state-of-the-art techniques based on sparse regularization not only offer similar accuracy, but are also much more stable to the given perturbations.

The instability phenomenon in DL for inverse problems is on the one hand similar, but on the other hand different to the better-known phenomenon of adversarial attacks in DL for classification problems [13]. In both classes of problem, such susceptibilities have potentially serious consequences. Indeed, there are at least two distinct areas where this issue may have significant impact. First, areas where DL techniques are designed to perform tasks hereto performed by humans. For example, automatic diagnosis in medicine, as has already been approved for commercial use in April 2018 by the US Food and Drug Administration (FDA) [19]. Similarly, a recent publication in Science [20], for example, warns about the potentially severe consequences in insurance fraud. Second, as noted above, areas where DL techniques replace well-established algorithms in the computational sciences. Instabilities therein may lead to incorrect scientific predictions, with serious downstream consequences.

¹To be specific, [44] is titled 'AI transforms image reconstruction' and features a new DL approach [50] which 'improves speed, accuracy and robustness of biomedical image reconstruction'.

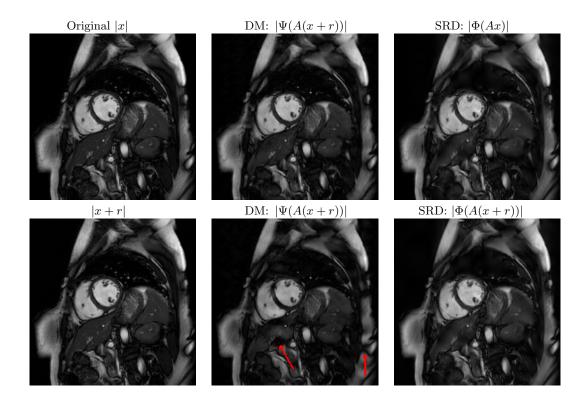


Figure 2: **(False Positives).** Perturbations cause the Deep MRI network (DM) $\Psi: \mathbb{C}^m \to \mathbb{C}^N$ [41] to reconstruct false positives (red arrows). Left: Original image x (top) and perturbed image x+r (bottom). Middle: DM reconstructs the image x from measurements y=Ax (top) and $\widetilde{y}=A(x+r)$ (bottom). Here $A\in\mathbb{C}^{m\times N}$ is a subsampled Fourier transform. Right: A sparse regularization decoder (SRD) $\Phi:\mathbb{C}^m\to\mathbb{C}^N$ reconstructs the image x from measurements y=Ax (top) and $\widetilde{y}=A(x+r)$ (bottom). This experiment is from [6].

With this in mind, the objective of this paper is to initiate a rigorous mathematical programme for investigating, understanding and eventually mitigating instabilities in DL for problems in computational science. Due to the impact it has had so far, our focus shall be on the case of inverse problems, and in particular, inverse problems in imaging. Our main contributions are outlined next.

1.2 Contributions

We consider the following finite-dimensional, underdetermined inverse problem:

Given measurements
$$y = Ax$$
, recover x . (1.1)

Here $A \in \mathbb{C}^{m \times N}$, m < N, is the sampling operator (or measurement matrix), $y \in \mathbb{C}^m$ is the vector of measurements and $x \in \mathbb{C}^N$ is the object to recover (typically a vectorized discrete image). To avoid pathological examples we assume that $\operatorname{rank}(A) \geqslant 1$. Our objective is to construct a reconstruction method; that is, a map taking the measurements y as input and giving (an approximation to) x as output.

Our focus in this work is on the underdetermined setting m < N, as is now common in practice. In many applications, the number of measurements is severely limited, due to time, cost, power or other constraints, hence one commonly faces the situation where $m \ll N$. To design a good reconstruction mapping one generally requires additional information on the images to be reconstructed. Sparse regularization techniques, described in more detail in §4.1, are typically untrained: they exploit the inherent sparsity of natural images in fixed transform domains (e.g. wavelets or discrete gradient). Over the last fifteen years, such techniques, supported by

the theory of *compressed sensing*, have become the state-of-the-art for many different image reconstruction tasks. On the other hand, more recent DL techniques, introduced formally in §2, are *data-driven*: they seek to implicitly learn a suitable image structure from a training database of existing images.

1.2.1 The instability phenomenon in DL for inverse problems

Fig. 1 demonstrates one type of instability for DL methods in inverse problems: namely, visually imperceptible perturbations in the input causing large effects on the output. There is also another type of instability, also introduced in [6]. Here, a small structural change in the image that is clearly visible and typically significant is highly distorted, or even removed altogether in the reconstructed image. If that image is used for diagnosis, such a perturbation can conceivably lead to a false positive or a false negative. Fig. 2 gives an example of this phenomenon: an imperceptible change in the image causes an artefact in the reconstruction that cannot easily be dismissed as unphysical. A more dramatic example is shown in Fig. 3. The network therein was trained to recover ellipses from Radon measurements. Yet, when the image is contaminated by a visible detail that is not part of the training set, it is completely washed out by the network. Notice that the addition of this detail also causes significant distortion in the recovery of the ellipse located directly above it. Since the inserted detail was not part of the training set, one may be inclined to dismiss its poor recovery by the network. However, ellipses were certainly in the training set, so the distortion seen in the recovery of this feature is troubling. This example also points towards the difficulty in understanding the generalization performance of trained NNs.

With these instabilities in mind, in our first contribution we seek to address the following question:

I. Why are DL methods for inverse problems typically unstable?

Our key insight into this question is kernel awareness. Loosely speaking, a reconstruction method lacks kernel awareness if it approximately recovers two vectors x and x' (from their respective measurements) whose the difference x'-x lies either in or close to the null space of A. Specifically, our main findings are:

- (i) Instabilities and false positives/negatives. DL can be unstable because typical training procedures lead to methods lacking kernel awareness. This causes both arbitrarily large Lipchitz constants and false positives and negatives. A false positive is a detail, for example a brain tumour in medical imaging, that is not present in the original image, but is present in the reconstructed image. Similarly, a false negative is a detail that is present in the original image, but is washed out by the reconstruction method.
- (ii) Instabilities are typically stable. The perturbations that lead to poor reconstructions and false positives/negatives do not belong to a set of Lebesgue measure zero. In fact there are sets, containing balls, of perturbations that will result in the method producing poor reconstructions. Hence, perturbing the perturbation that creates the unwanted recovery will also yield a poor result. In other words, the instabilities are typically stable.
- (iii) **Instabilities are not rare events.** For much the same reason, damaging perturbations are not rare events. Under reasonable random noise models, they occur with nonzero probabilities.

Finding (iii) is highlighted in Fig. 4. While the perturbations shown in Figs. 1 and 2 are designed specifically to cause severe artefacts, as Fig. 4 demonstrates, mean-zero, random perturbations of the measurements do, with nonzero probability, lead to similar effects.

Our final two findings in this section examine the ubiquity of this phenomenon:

(iv) The instability phenomenon occurs regardless of the architecture. The instability phenomenon, as described in (i), is not related to a particular network architecture.

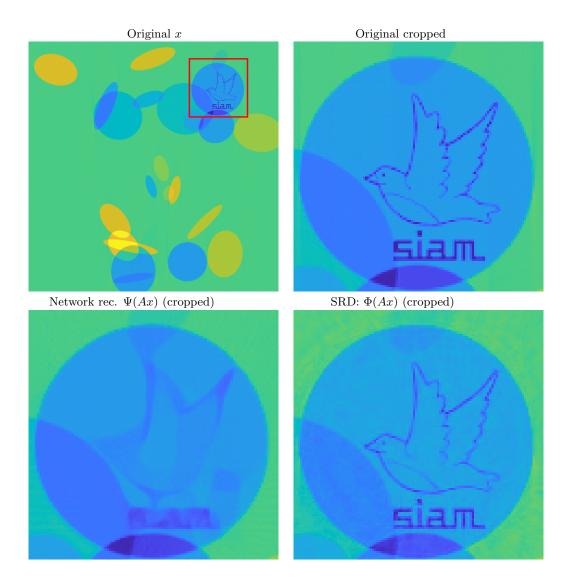


Figure 3: **(False negatives).** The FBPConvNet map $\Psi: \mathbb{R}^m \to \mathbb{R}^N$ [28] is trained to recover images comprised of ellipses from a Radon sampling operator $A \in \mathbb{R}^{m \times N}$. Top left: The image x containing a bird and the SIAM logo. This is a feature the network has not seen. Top right: Cropped original image. Lower left: The cropped FBPConvNet reconstruction from measurements Ax. Lower right: The cropped reconstruction of x from measurements Ax using a sparse regularization decoder (SRD) $\Phi: \mathbb{R}^m \to \mathbb{R}^N$. See Section 9 for further details.

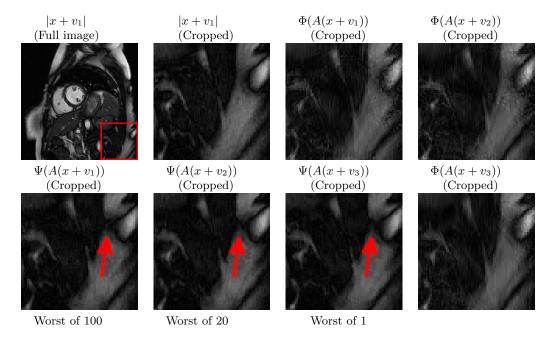


Figure 4: (False positives and random noise). The DM network of Fig. 2 is unstable with respect to Gaussian noise. We compute 100 Gaussian noise vectors $w_j \in \mathbb{C}^N$ and pick (using the eyeball metric) the w_j for which $\Psi(A(x+w_j))$ gives the larges artefact, and label it v_1 . Repeating the experiment with 20 (and 1) new noise vectors yields a perturbation v_2 (and v_3). Ψ introduces a false dark area, indicated by the red arrows. The SRD of Fig. 2 is applied to the same measurements $A(x+v_j)$, j=1,2,3. See Section 9 for further details.

Namely, it is the phenomenon of lack of kernel awareness, as described above, that is the main contributing factor.

(v) The instability phenomenon is not easy to cure. The presence of instabilities in existing DL methods for inverse problems in imaging naturally leads one to consider how they may be overcome. Our main result sheds light on this issue, in particular, it allows one to predict which strategies are unlikely to work. It also highlights that avoiding instabilities is not straightforward, due to mild conditions established in the theorem which cause them to appear.

These findings are discussed in §3, and specifically, Theorem 3.1.

1.2.2 Optimal reconstruction maps for inverse problems

As discussed above, modern techniques for solving (1.1) exploit additional information about the images to be reconstructed. Hence, mathematically, it is now convenient to describe the inverse problem in terms of a triple $\{A, \mathcal{M}_1, \mathcal{M}_2\}$, given by

a sampling map
$$A \in \mathbb{C}^{m \times N}$$
, where $m < N$,
a domain $\mathcal{M}_1 \subset \mathbb{C}^N$, where (\mathcal{M}_1, d_1) is a metric space,
the range $\mathcal{M}_2 = A(\mathcal{M}_1) \subset \mathbb{C}^m$, where (\mathcal{M}_2, d_2) is also a metric space. (1.2)

The inverse problem is now as follows:

Given measurements
$$y = Ax$$
 of $x \in \mathcal{M}_1$, recover x . (1.3)

In methods such as sparse regularization, the domain \mathcal{M}_1 is chosen explicitly as a suitable a priori model for the class of images, typically by imposing (approximate) sparsity in some transform, and the reconstruction method is designed specifically with \mathcal{M}_1 in mind. Conversely,

in data-driven methods such as DL, one does not choose \mathcal{M}_1 explicitly, but rather has access to samples of it via a training set

$$\mathcal{T} := \{ (y^j, x^j) : y^j = Ax^j, x^j \in \mathcal{M}_1, j = 1, \dots, K \}.$$
(1.4)

In either case, the objective is to construct a reconstruction mapping

$$\Phi: \mathbb{C}^m \to \mathbb{C}^N$$
,

tailored to approximate $x \in \mathcal{M}_1$ from its measurements y = Ax as accurately as possible. This raises the following basic question:

II. Given a triple $\{A, \mathcal{M}_1, \mathcal{M}_2\}$ what is the optimal reconstruction map $\Phi : \mathbb{C}^m \to \mathbb{C}^N$?

In order to address this question, in §5 we introduce the concept of an *optimal map*. Our main finding, described in Theorem 5.2, is as follows:

(vi) Training may not yield an optimal map. The mapping obtained through the training process may not be optimal, or even approximately optimal. Hence, there is no guarantee that DL will perform better than sparse regularization. Moreover, the mapping sought by the training procedure may not actually exist. In particular, one may be in the situation that there does not exist a NN that can, given the measurements, obtain small training and test errors.

To elaborate on this finding, as discussed further in $\S 2$, recall that the objective of training is to construct a mapping Ψ which performs well on the a set \mathcal{U} ,

$$\|\Phi(y) - x\| \leqslant \delta, \quad \forall (y, x) \in \mathcal{U},$$
 (1.5)

where \mathcal{U} contains at least the training set (1.4), but also, ideally, the so-called *test set*. The above finding states firstly that any such a map may not be an optimal map for $\{A, \mathcal{M}_1, \mathcal{M}_2\}$, and secondly, that it may not even be possible to achieve (1.5) in the first place. We remark in passing that this phenomenon is not the same as the well-known concept of *overfitting* in machine learning. It relies crucially on the fact that A is noninvertible. See §5.2 for further details.

It is well known that the practical performance of NNs can often be enhanced by adding a regularization term to the objective function used in the training process (see §2). This functional involves a regularization *parameter*, used to balance the data fidelity and regularization terms, which requires tuning so as to get good performance. Our next finding, described in §7, states that this is a highly delicate issue:

(vii) Performance can be stable or unstable with respect to the regularization parameter. There are many cases where a suitable choice of the regularization parameter in DL yields an optimal map. Yet, for some problems, if a new element is added to the the training set or replaces an existing element, then there is no choice of regularization parameter which yields an optimal map. Conversely, for other problem, one may add arbitrarily many elements to the training step and obtain an optimal map using the same regularization parameter. Hence, the dependence of the performance of DL on the training parameter can be completely brittle or completely stable depending on the problem.

1.2.3 State-of-the-art performance for inverse problems

As discussed further in §4, reconstruction techniques for inverse problems in imaging based on sparse regularization possess, under suitable conditions, guarantees on their stability and accuracy through the theory of compressed sensing. In the third component of this work, we consider the following question:

Our main findings, elaborated in §6, are as follows:

- (viii) **DL** may outperform sparse regularization at the cost of instability. There are cases where DL attains lower errors than sparse regularization, but in doing so it is unstable. Moreover, the bigger the improvement the more unstable it becomes.
- (ix) Two additional training samples may destabilize DL. There are large classes of Fourier inverse problems for which DL and sparse regularization will both give optimal maps. However, adding two elements to the training set (that can be arbitrarily close to the existing training data), destabilizes DL and renders it non-optimal.

Note that a Fourier inverse problem corresponds to the case where A is a subsampled discrete Fourier transform. This setup occurs notably in applications such as MRI.

1.3 Related work

Sparse regularization rose to prominence in the imaging community with the introduction of compressed sensing in the 2000's by Candès, Romberg & Tao [11] and Donoho [16]. See [21] for an in-depth treatment. The theory of compressed sensing has been much lauded and it lies at the heart of much of this paper. As a counterpoint to findings (i)–(v) above, we observe in §4 that conditions such as the Restricted Isometry Property [12] and robust Null Space Property [21] which ensure stable and accurate recovery with compressed sensing, are precisely statements about its kernel awareness. Our work on optimal maps (see (vi)) is inspired by the seminal work of Cohen, Dahmen & DeVore [14] on optimal decoders for compressed sensing, and in particular, the notion of instance optimality.

Instabilities in DL for classification problems were first discovered in [46]. A significant development was the *DeepFool* package of Moosavi–Dezfooli, Fawzi & Frossard [33], which was followed by the construction of so-called *universal adversarial perturbations* [32]. The construction of and mitigation against *adversarial attacks* is now an active area of research. To the best of our knowledge, [6] and [27] were the first works to demonstrate the instability phenomenon for inverse problems in imaging.

The work of Jin, McCann, Froustey & Unser [28] was influential in highlighting the promise of DL for inverse problems in imaging. This is now a rapidly evolving area of research, which we will not attempt to summarize. See [31], [5] and [7] for overviews of current techniques. Note that sparse regularization has been used as the basis for some DL technique, e.g. by using DL to recover the parts of an image that sparse regularization cannot such as in [10], or by designing NN architectures through the process of *unravelling* an optimization algorithm (see, e.g., [7]).

Finally, let us note that this work is considers deterministic approaches to inverse problems. For an in-depth treatment of Bayesian approaches, see, for instance the work of Stuart [45] and Dashti & Stuart [15], as well as references therein.

1.4 Summary and outlook

The purpose this paper is to initiate a programme into the rigorous foundations of NNs and DL from the dual pillars of numerical analysis, accuracy and stability. While much of the focus on DL in the machine learning community has been on discrete problems such as classification, this paper aims to highlight both the challenges and potential when applying DL to continuous problems in computational mathematics. As shown in $\S 5.2$, such problems can be fundamentally different to the classification problem. Due to both their ubiquity in the computational sciences and the recent activity on data-driven approaches for them, we have chosen to focus on inverse problems. Yet we expect both questions I–III and our findings (i)–(ix) to be relevant to other problems in computational mathematics.

For inverse problems, the conclusions from our findings are decidedly mixed: current approaches to training cannot ensure stable methods; even if they do, the resulting methods may

not offer state-of-the-art performance; regularization strategies may not fix these issues. Furthermore, these are not rare events, able to be dismissed by all, but a small group of theoreticians (recall Fig. 4). Should one therefore give up on the DL approach to inverse problems? Of course not. The rich approximation theory – dating back to the classical $Universal\ Approximation\ Theorem$ (see, e.g., [35]) but including many recent advances such as [42, 49] – says that NNs have the potential to give rise to powerful methods for inverse problems in imaging. Our hope is that these findings, in particular, the crucial role that $kernel\ awareness$ plays, spur new research into devising better ways to design and train stable and accurate DL algorithms.

1.5 Outline

The outline of the remainder of this paper is as follows. In §2 we give a short introduction to DL. In §3 we develop the instability phenomenon, and in §4 we consider kernel awareness in the context of sparse regularization. We develop optimal maps in §5, and in §6 we compare DL and sparse regularization from the perspective of stability and optimality. In §7 we discuss regularization for DL. Finally, in §8 we have gathered proofs of the main results and §9 implementation details for the numerical experiments can be found. Code and data accompanying this paper is available at https://github.com/vegarant/troub_ker.

2 Deep learning overview

We commence this paper with a very short introduction to DL. See, e.g. [22, 25], for more comprehensive treatments.

The objective of DL is to construct a neural network that approximates a map $f: \mathbb{R}^d \to \mathbb{R}^p$ from samples, i.e. pairs (y, f(y)). There are many ways to define a neural network in practice. In this paper, we consider the vanilla case of feedforward neural networks, although we note in passing that many of our results also apply to more exotic setups. Recall that an L-layer feedforward neural network is a function $\Psi: \mathbb{R}^d \to \mathbb{R}^p$ of the form

$$\Psi(x) = V_L(\rho(V_{L-1}(\rho(\ldots \rho(V_1(x)))))), \quad x \in \mathbb{R}^d,$$

where each $V_j: \mathbb{R}^{n_{j-1}} \to \mathbb{R}^{n_j}$ is an affine map

$$V_j x = W_j x + b_j, \qquad W_j \in \mathbb{R}^{n_j \times n_{j-1}}, \quad b_j \in \mathbb{R}^{n_j},$$

 $\rho: \mathbb{R} \to \mathbb{R}$ is a non-linear function, $\rho(x) = (\rho(x_i))$ for $x = (x_i)$, and $n_0 = d$, $n_L = p$. The W_j 's are referred to as weights and the b_j 's as biases. The number L is the depth of the network, and n_l is the width of its lth layer. The function ρ is the activation function. Typical choices for ρ are the Rectified Linear Unit (ReLU), defined by $\rho(x) = \max\{0, x\}$, or the sigmoid, defined by $\rho(x) = \frac{1}{1+e^{-x}}$.

The architecture of a neural network refers to choice of the depth L, widths n_1, \ldots, n_{L-1} and activation function ρ . Write $\mathbf{n} = (n_0, n_1, \ldots, n_L)$. We denote the class of neural networks with a given architecture as $\mathcal{N}\mathcal{N}$, or $\mathcal{N}\mathcal{N}_{\mathbf{n}}$ when we wish to make the dependence on \mathbf{n} explicit. Note that $\mathcal{N}\mathcal{N}_{\mathbf{n}}$ is parametrized by the weights and biases $\{W_j, b_j\}_{j=1}^L$.

In inverse problems, our goal is to construct a mapping that takes as input the measurements y = Ax of some unknown image x, and returns x (or some approximation to it) as output. In DL for inverse problems (see, e.g., [7, 28, 31]) this is achieved using a training set

$$\mathcal{T} = \{ (y^j, x^j) : y^j = Ax^j, \ j = 1, \dots, K \},\$$

consisting of pairs of the form (Ax, x), where x is a training image and y = Ax are its measurements.

Remark 2.1. In inverse problems (typically MRI), it is common to deal with complex input $y \in \mathbb{C}^m$ and output $x \in \mathbb{C}^N$. A standard way is to associate $y \in \mathbb{C}^m$ with a vector $y' \in \mathbb{R}^{2m}$ consisting of the real and imaginary parts of y, then apply a real-valued network $\Psi : \mathbb{R}^{2m} \to \mathbb{R}^{2N}$. Similarly $x' = \Psi(y')$ is associated with a complex image $x \in \mathbb{C}^N$. We assume the complex case is treated in this way throughout this paper. Henceforth we simply write $\Psi : \mathbb{C}^m \to \mathbb{C}^N$ for a network taking complex inputs and outputs, with the assumption that it has this form.

Having fixed an architecture, i.e. a class \mathcal{NN} , training a neural network is the process of computing an approximation $\Psi \in \mathcal{NN}$ from the data \mathcal{T} . This is typically achieved by computing a minimiser of a certain optimization problem

$$\Psi \in \underset{\widetilde{\Psi} \in \mathcal{NN}}{\operatorname{argmin}} \frac{1}{K} \sum_{j=1}^{K} \operatorname{Cost}(\widetilde{\Psi}, Ax^{j}, x^{j}), \tag{2.1}$$

where Cost is an appropriate cost function. A popular choice is the ℓ^2 -loss, given by

$$\mathrm{Cost}(\widetilde{\Psi},Ax^j,x^j) = \frac{1}{2} \|x^j - \widetilde{\Psi}\left(Ax^j\right)\|_{\ell^2}^2.$$

However, there is a myriad of other possibilities.

The optimization problem is generally nonconvex (2.1), and devising effective numerical methods for computing minima is a substantial topic in its own right. This topic is not the concern of this paper. Instead, we simply observe that the result of a such a procedure is, in general, a neural network $\Psi \in \mathcal{NN}$ with, at the very least, small training error. That is,

$$\|\Psi(y) - x\| \le \delta, \qquad \forall (y, x) \in \mathcal{T},$$
 (2.2)

for some $\delta > 0$ and some norm $\|\cdot\|$. In practice, δ will depend on the cost function and the algorithm used for approximately solving (2.1).

An common issue in DL is *overfitting*. This refers to a network Ψ with small training error, but poor *generalization*: that is, the trained network performs poorly on new images x that are not in the training set. Typically, generalization performance is measured using a second set of data, the *test set*, not used in the training process. The corresponding error is the *test error*.

Regularization is standard means to attempt to cure the tendency of trained NNs to overfit. Instead of (2.1), one attempts to compute

$$\Psi \in \underset{\widetilde{\Psi} \in \mathcal{NN}}{\operatorname{argmin}} \frac{1}{K} \sum_{i=1}^{K} \frac{1}{2} \|x_i - \widetilde{\Psi}(Ax_i)\|_{\ell^2}^2 + \lambda J(\widetilde{\Psi}), \tag{2.3}$$

where $\lambda \geqslant 0$ and $J : \mathcal{NN} \to \mathbb{R}$ is a regularization function. Often J may be chosen to penalize large weight matrices.

3 Lack of kernel awareness: why DL methods may lead to unstable methods for inverse problems

Figs. 1–4 highlight the instability of several existing DL algorithms for inverse problems. Of course, unstable algorithms can arise for many reasons – for instance, a poor numerical implementation. However, in this setting, there is a more fundamental explanation involving the null space of A. Specifically, a stable recovery algorithm must not commit the *cardinal sin* of recovering images whose difference lies either in or close to $\mathcal{N}(A)$. If it does, then, as shown next, it is necessarily unstable.

3.1 Instabilities, false positives and false negatives

Theorem 3.1 below states that any method committing the above sin admits a lower bound on its local Lipschitz constant around a certain vector that is reconstructed with small error. Under certain conditions, which are easily satisfied in the case of DL using standard training, this bound can become arbitrarily large. Hence, a small perturbation can cause large artefacts in the reconstructed image. In addition, Theorem 3.1 also guarantees the existence of both false positives and negatives under the same conditions. Thus, it also explains the vulnerability of DL to small structural changes in the input.

Striving for generality, in the following result we endow \mathbb{C}^N and \mathbb{C}^m with metrics d_1 and d_2 respectively. Let $\Phi: \mathbb{C}^m \to \mathbb{C}^N$ be a reconstruction map. Then we define the ε -Lipschitz constant of Φ at $y \in \mathbb{C}^m$ as

$$L^{\varepsilon}(\Phi, y) = \sup_{0 < d_2(z, y) \leqslant \varepsilon} \frac{d_1(\Phi(z), \Phi(y))}{d_2(z, y)}.$$

Theorem 3.1 (The cardinal sin: recovering elements close to $\mathcal{N}(A)$). Let $A: \mathbb{C}^N \to \mathbb{C}^m$ be a linear map, d_1 and d_2 be metrics on \mathbb{C}^N and \mathbb{C}^m respectively, and $\Psi: \mathbb{C}^m \to \mathbb{C}^N$ be continuous. Suppose that there exist $x, x' \in \mathbb{C}^N$ and $\eta > 0$ such that

$$d_1(\Psi(Ax), x) < \eta, \quad d_1(\Psi(Ax'), x') < \eta,$$
 (3.1)

and

$$d_2(Ax, Ax') \leqslant \eta. \tag{3.2}$$

Then the following hold:

(i) (Instability). There is a closed non-empty ball $\mathcal{B} \subset \mathbb{C}^m$ centred at y = Ax such that the local ε -Lipschitz constant at any $\widetilde{y} \in \mathcal{B}$ is bounded from below:

$$L^{\varepsilon}(\Psi, \widetilde{y}) \geqslant \frac{1}{\varepsilon} (d_1(x, x') - 2\eta), \quad \varepsilon \geqslant \eta.$$
 (3.3)

(ii) (False positives). Suppose that the metrics are translation invariant. Then there exists $z \in \mathbb{C}^N$, $e \in \mathbb{C}^m$, with $d_1(0,z) \geqslant d_1(x,x')$, $d_2(0,e) \leqslant \eta$, and closed non-empty balls \mathcal{B}_x , \mathcal{B}_e and \mathcal{B}_z centred at x, e and z respectively such that

$$d_1(\Psi(A\widetilde{x} + \widetilde{e}), \widetilde{x} + \widetilde{z}) \leqslant \eta, \quad \forall \widetilde{x} \in \mathcal{B}_x, \ \widetilde{e} \in \mathcal{B}_e, \ \widetilde{z} \in \mathcal{B}_z. \tag{3.4}$$

(iii) (False negatives). Suppose that the metrics are translation invariant. Then there exists $z \in \mathbb{C}^N$, $e \in \mathbb{C}^m$, with $d_1(0,z) \geqslant d_1(x,x')$, $d_2(0,e_1) \leqslant \eta$, and closed non-empty balls \mathcal{B}_x , \mathcal{B}_e and \mathcal{B}_z centred at x, e and z respectively such that

$$d_1(\Psi(A(\widetilde{x}+\widetilde{z})+\widetilde{e}),\widetilde{x}) \leqslant \eta, \quad \forall \widetilde{x} \in \mathcal{B}_x, \ \widetilde{e} \in \mathcal{B}_e, \ \widetilde{z} \in \mathcal{B}_z. \tag{3.5}$$

As this result shows, not only is it possible to develop instabilities, it is alarmingly easy to do so. Simply performing well on two vectors whose difference lies close to the null space $\mathcal{N}(A)$ is enough. This situation is not rare. For instance in MRI, the matrix A is a subsampled discrete Fourier transform. Typical MRI images may be of size $N=256\times256=65536$, and it is common to subsample by a factor of 25%, giving m=16384. Hence $\mathcal{N}(A)$ has dimension N-m=49152. Thus, given a training set of typical MRI images, the large dimension of the null space will typically offer a multitude of ways to achieve (3.2).

This theorem also illustrates the ease with which a trained method can produce false positives and negatives. This is shown in Fig. 5. It shows the recovery of an image x and an image x', which is identical to x except for a small structural change representing a tumour. The tumour, x'-x, is localized in space and its Fourier transform is approximately zero at low frequencies. The measurement matrix A is a subsampled discrete Fourier transform, and takes mainly low frequency measurements (as is typical in practice). Hence x'-x lies approximately in $\mathcal{N}(A)$. An NN Ψ is trained to recover both x' and x well from their respective measurements. But, as the theorem implies, and this figure confirms, the result is false positives (a small perturbation of y = Ax causes the presence of a tumour in the reconstruction) and false negatives (a small perturbation of y' = Ax' removes the tumour in the reconstruction).

Remark 3.2 (Metrics in Theorem 3.1). Theorem 3.1 is deliberately phrased using metrics in order to demonstrate that regardless of the distance function, the instability phenomenon will occur. Norms are arguably commonly used as error functions in inverse problems, however, it is well known that errors measured in distance functions induced by a norm may not always represent what the human eye interprets as either large or small errors. However, if the 'eyeball metric' could be represented by a mathematical metric, then Theorem 3.1 holds for it. Throughout the paper we state theorems for general metrics whenever the proof allows for such generality.

Theorem 3.1 is deliberately formulated with the very weak conditions (3.1) and (3.2) so as to demonstrate the ease with which a method can become unstable. However, although the instabilities always happen in a ball, Theorem 3.1 does not say anything about how big these balls are. For this one needs slightly stronger assumptions on the of the mapping Ψ :

Corollary 3.3 (Instabilities in larger sets). Let $A \in \mathbb{C}^{m \times N}$, $m \leq N$, be full rank and let $\Psi : \mathbb{C}^m \to \mathbb{C}^N$ be continuous. Consider \mathbb{C}^m and \mathbb{C}^N equipped with their respective ℓ^2 -norms. Suppose that there exists $x \in \mathcal{N}(A)^{\perp}$, $x' \in \mathbb{C}^N$ and $r_1 > 0$ such that

$$\|x' - \Psi(Ax')\|_{\ell^2} < \eta, \quad \|z - \Psi(Az)\|_{\ell^2} < \eta, \quad \|A(x' - z)\|_{\ell^2} \leqslant \eta, \quad \forall z \in \mathcal{B}_x,$$

where \mathcal{B}_x is the open ball of radius r_1 centred at x. Then there exists a closed ball $\mathcal{B}_y \subset \mathbb{C}^m$ of radius $r_2 \geqslant \sigma_{\min}(A)r_1$ centred at y = Ax such that the following holds. For every $\widetilde{y} \in \mathcal{B}_y$, the local ε -Lipschitz constant satisfies

$$L^{\varepsilon}(\Psi, \widetilde{y}) \geqslant \frac{1}{\varepsilon} \left(\operatorname{dist}(x', \mathcal{B}_x) - 2\eta \right), \quad \varepsilon \geqslant \eta,$$

where $\operatorname{dist}(x', \mathcal{B}_x) = \inf\{\|x' - z\|_{\ell^2} : z \in \mathcal{B}_1\}$ and $\sigma_{\min}(A)$ is the smallest singular value of A.

Remark 3.4 (The instabilities are stable). As Theorem 3.1 demonstrates, even in the most general case, there is always a ball around the perturbation that causes the instability such that any perturbation within this ball also yields an instability. In particular, the instability is stable. With stronger assumptions such as in Corollary 3.3, the size of the balls causing instabilities can be quantified. This effect is illustrated in Fig. 6, based on the example shown previously in Fig. 1.

Note that the conditions (3.1) in Theorem 3.1, namely,

$$d_1(\Psi(Ax), x) < \eta, \quad d_1(\Psi(Ax'), x') < \eta, \quad \eta > 0$$

are exactly what is expected when training a neural network. As discussed previously, training typically yields as small training error, implying

$$\|\Psi(Ax) - x\| \le \delta, \quad \forall (Ax, x) \in \mathcal{T},$$

for some $\delta > 0$, some norm $\|\cdot\|$ and some training set \mathcal{T} . Theorem 3.1 simply says that if the training set has at least two elements (Ax, x) and (Ax', x') for which $\|A(x - x')\|$ is small then instabilities necessarily occur.

Moreover, the assumption in Corollary 3.3 that $A \in \mathbb{C}^{m \times N}$ with $m \leqslant N$ be full rank is not rare. In MRI, for example, the measurement matrix takes the form $A = P_{\Omega}F$ where F is the discrete Fourier transform, $\Omega \subseteq \{1, \ldots, N\}$, $|\Omega| = m$ is a sampling set (a set of indices corresponding to the frequencies measured) and P_{Ω} is a matrix that selects the indices according to Ω . Observe that $\sigma_{\min}(A) = 1$ in this case, since F is unitary and Ω is generally assumed to not contain any repeats.

Note that throughout the paper, we consider P_{Ω} to be either an $m \times N$ matrix, so that A is $m \times N$, or an $N \times N$ matrix, i.e. the orthogonal projection onto span $\{e_j : j \in \Omega\}$, where $\{e_j\}$ is the canonical basis of \mathbb{C}^N . The precise meaning will be clear from the context.

Remark 3.5 (Stability implies a universal barrier on performance). Theorem 3.1 describes a basic mechanism that causes instabilities in learning for inverse problems. Indeed, for a method to be stable there must be some awareness of the kernel of the sampling operator A. The theorem reveals that there is a barrier restricting how well any reconstruction method can perform. Indeed, if one recovers two vectors well whose difference is close to the kernel (precisely the regime that one should not be able to recover from the data) the instability phenomenon will occur and there are inputs that will give false positives and false negatives.

Remark 3.6 (The Lipschitz constant can be arbitrary large). It is clear from (3.3) that the ε -Lipschitz constant can become arbitrarily large. Indeed, even when we have $\dim(\mathcal{N}(A)) = 1$, there will be examples of training sets where η can be arbitrarily small while $d_1(x,x')$ stays bounded from below, hence, $L^{\varepsilon}(\Phi,\widetilde{y})$ can be arbitrary large. Note that any subsampled inverse problem in which A is an $m \times N$ matrix and m < N obviously yields $\dim(\mathcal{N}(A)) \geqslant 1$.

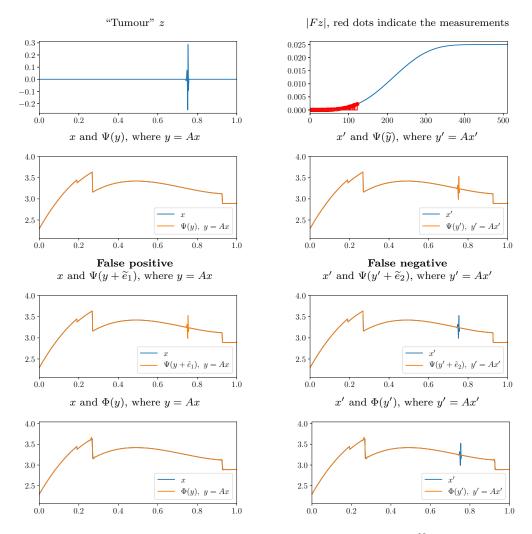


Figure 5: (False positives/negatives). A neural network $\Psi\colon \mathbb{C}^m\to\mathbb{R}^N$ is trained to approximately recover x and x'=x+z from measurements y=Ax and y'=Ax' respectively, where $\|z\|_{\ell^2}=0.4$ and $\eta=\|A\widetilde{z}\|_{\ell^2}=0.005566$ and $A=P_\Omega F$, where $F\in\mathbb{C}^{N\times N}$ is a discrete Fourier transform. First row: Tumour z and the magnitude of each component in $|F\widetilde{z}|$. The red dots indicate the samples corresponding to Ω . Second row: Original image x and the image x'=x+z containing the tumor z. Third row: Perturbations \widetilde{e}_1 and \widetilde{e}_2 are added to the measurements y=Ax and y'=Ax' respectively, which induce the network to produce a false positive and a false negative respectively. Fourth row: The recovery of the images from the measurements y and y' by a sparse regularization decoder $\Phi\colon\mathbb{C}^m\to\mathbb{R}^N$. Observe that Φ does not recover z as it lies close to $\mathcal{N}(A)$.

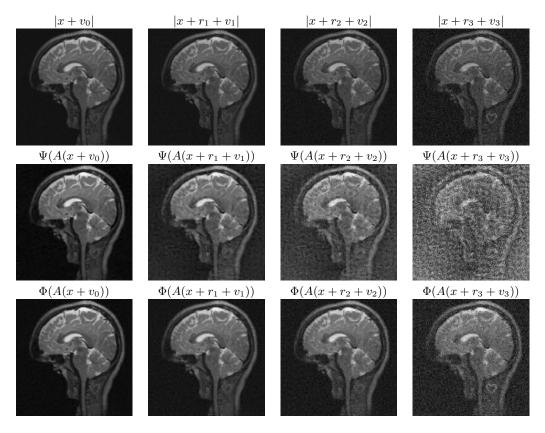


Figure 6: (The instabilities are stable). Adding a random perturbation to a 'bad' perturbation yields a 'bad' perturbation. Based on the experiment in Fig. 1, new perturbations $||v_0||_2 = ||v_1||_2 < ||v_2||_2 < ||v_3||_2$ are constructed with magnitude $||v_i||_2/||r_i||_2 = 1/2$ for i = 1, 2, 3. The $v_i s$ are constructed as $v_i = A^* e_i$ where the real and imaginary components of $e_i \in \mathbb{C}^m$ are drawn independently from a normal distribution $\mathcal{N}(0,1)$ and then rescaled to get the desired norm. Top row: The magnitude of the sampled images $|x + r_i + v_i|$, for i = 0, 1, 2, 3, with $r_0 = 0$. Middle row: The reconstruction obtained by the AUTOMAP network $\Psi : \mathbb{C}^m \to \mathbb{C}^N$, given measurements $A(x + r_i + v_i)$, i = 0, 1, 2, 3. Third row: Reconstruction from $A(x + r_i + v_i)$, i = 0, 1, 2, 3 using a sparse regularization decoder $\Phi : \mathbb{C}^m \to \mathbb{C}^N$ (see Section 9 for details).

3.2 Do bad perturbations occur in practice?

Theorem 3.1 is a worst-case bound. It asserts the existence of a damaging perturbation of the measurements y = Ax which may cause large artefact and/or false positives and false negatives. But do such perturbations actually occur in practice? Or are they 'rare events', of interest only to the theorist, but not the practitioner? In the following we will discuss these questions.

In practice, perturbations often arise as random noise on the measurements, i.e. realizations of a mean-zero random variable $e_{\text{pert}}:\Omega\to E$, where $(\Omega,\mathcal{F},\mathbb{P})$ is a probability space and $E=\mathbb{C}^m$ equipped with the metric d_2 (induced by some norm) and the Borel measure. The fact that (3.3), (3.4) and (3.5) in Theorem 3.1 hold for all perturbations within a ball implies lower bounds on the probability of having 'bad' perturbations realized by e_{pert} . In particular, given the assumptions of Theorem 3.1, and with mild assumptions on e_{pert} , we have

$$\mathbb{P}\left(L^{\varepsilon}(\Psi, y + e_{\text{pert}}) \geqslant \frac{1}{\varepsilon} \left(d_1(x, x') - 2\eta\right)\right) \geqslant c_1 > 0, \tag{3.6}$$

where y = Ax, as well as

$$\mathbb{P}\left(d_1(\Psi(y+e_{\text{pert}}), x+z) \leqslant \eta\right) \geqslant c_2 > 0,\tag{3.7}$$

and

$$\mathbb{P}\left(d_1(\Psi(A(x+z) + e_{\text{pert}}), x) \leqslant \eta\right) \geqslant c_3 > 0. \tag{3.8}$$

Here z is the false positive or false negative obtained in Theorem 3.1 and c_1, c_2, c_3 are constants depending on e_{pert} and the sizes of the balls guaranteed by Theorem 3.1. These observations are illustrated numerically in Fig. 4. Note that c_1, c_2, c_3 above can be estimated more explicitly given further assumptions on e_{pert} and bounds on the sizes of these balls; see Corollary 3.3 for the latter.

This discussion pertains to *generic* noise, i.e. realizations of mean-zero, random variables. Yet in many applications, the measurement are corrupted not only by generic noise, but also by other phenomena. This is the case for instance in MRI, where factors such as small patient motion or small anatomic differences cause specific corruptions in the measurements. In such settings, a more suitable model of the random perturbation is

$$e_{\mathrm{pert}} = e_{\mathrm{pert}}^1 + e_{\mathrm{pert}}^2 : \Omega \to E,$$

where e_{pert}^1 is a random variable that accounts for the non-generic part of the perturbation and e_{pert}^2 is a mean-zero random variable accounting for the generic part.

While it is typically straightforward in applications to identify a reasonable model for $e_{\rm pert}^2$ (e.g. Gaussian, Poisson, etc), it is usually much less straightforward to model $e_{\rm pert}^1$. This, of course, is one motivation for establishing stability guarantees that address worst case perturbations. Conversely, as discussed in Remark 3.4 the instabilities established by Theorem 3.1 are stable. Hence, if a realization of $e_{\rm pert}^1$ results in a damaging perturbation, then (3.7) can be replaced by

$$\mathbb{P}\left(d_1(\Psi(y + e_{\text{pert}}), x + z) \leqslant \eta \mid e_{\text{pert}}^1 = e\right) \geqslant 1 - \varepsilon,\tag{3.9}$$

and similarly for (3.6) and (3.8). Since ε will typically be small, this is represents a high probability event: generic noise added to a damaging perturbation cannot counteract the damage. Such a phenonenon is shown empirically in Fig. 6.

This argument also demonstrates how difficult it may be to guarantee robustness to physical perturbations. In view of (3.9), doing so would involve ensuring that e_{pert}^1 rarely gives damaging perturbations and that $1 - \varepsilon$ is small. Since e_{pert}^1 is difficult to model, this may be practically impossible.

3.3 The instability phenomenon is not easy to remedy

Having established the presence of instabilities, the next question to ask is: how might one make DL more robust? There are many strategies one might try, yet proceeding in an ad-hoc fashion is both time- and resource-consuming. By establishing processes which lead to instabilities, Theorem 3.1 is a useful tool for excluding approaches that are unlikely to succeed in preventing instabilities. We now highlight three such strategies. The key idea is that any remedy which does not enforce some sort of kernel awareness will remain susceptible to instabilities.

Enforcing consistency

Consistency of the reconstruction with the measured data is often desirable in practice, and many emerging DL strategies for inverse problems seek to enforce this property [24, 26]. However, this does not prohibit instabilities. Indeed, let $\Psi: \mathbb{C}^m \to \mathbb{C}^N$ be an arbitrary reconstruction map (Ψ need not be a NN). Then consistency of Ψ , i.e.

$$A\Psi(y) = y, \quad \forall y = Ax, \, x \in \mathcal{M}_1.$$
 (3.10)

does nothing to help one avoid the conclusions of Theorem 3.1. Indeed, the corresponding conditions

$$d_1(\Psi(Ax), x) < \eta, \quad d_1(\Psi(Ax'), x') < \eta, \quad d_2(Ax, Ax') \le \eta,$$
 (3.11)

pertain to the quality of Ψ as an approximation, and are unrelated to its consistency. In fact, if $\Psi(Ax) = x$ and $\Psi(Ax') = x'$, i.e. Ψ recovers x and x' perfectly, then clearly Ψ is also consistent for x and x'.

A rather general approach to approximate consistency, as suggested in [24, 26], is to consider a set $\mathcal{S} \subset \mathbb{C}^N$ which either contains the images of interest, or approximates these images well. Then, one defines the reconstruction mapping $\Phi: \mathbb{C}^m \to \mathbb{C}^N$ as

$$y \mapsto \Phi(y) \in \underset{\widetilde{x} \in \mathcal{S}}{\operatorname{argmin}} \frac{1}{2} ||A\widetilde{x} - y||_{\ell^2}^2.$$
 (3.12)

However, if x and x' satisfy $d_2(Ax, Ax') \leq \eta$ and additionally, $x, x' \in \mathcal{S}$, then (3.11) still holds, and thus instabilities still occur.

Training with random sampling patterns

As noted, in applications such as MRI, the measurement matrix takes the form $A = P_{\Omega}F$, where $\Omega \subseteq \{1,\ldots,N\}$, $|\Omega| = m$ is the set of frequencies sampled and P_{Ω} is the projection onto these indices. Another approach to improve the robustness of DL, suggested in [43], is to train with many different sampling patterns at once. Specifically, one now considers the reconstruction map Ψ as a mapping from $\mathbb{C}^N \to \mathbb{C}^N$ and the measurement vector y as an element of \mathbb{C}^N , where the components in y that correspond to the unsampled indices are set to zero. The mapping Ψ is then found by training on the data $\{(y^{ji}, x^j) : j = 1, \ldots, K, i = 1, \ldots, L\}$, where $y^{ji} = P_{\Omega_i} x_j$ and Ω_i is the ith sampling pattern. For instance, one may define

$$\Psi \in \operatorname*{argmin}_{\widetilde{\Psi} \in \mathcal{NN}} \frac{1}{KL} \sum_{j=1,i=1}^{K,L} \frac{1}{2} \|x^i - \widetilde{\Psi} \left(P_{\Omega_j} F x^i\right)\|_{\ell^2}^2. \tag{3.13}$$

Once Ψ is trained, it is used to reconstruct an image x from measurements $y = P_{\Omega}F$ acquired from a given sampling pattern Ω of size $|\Omega| = m$ (which may or may not be equal to Ω_i for some i). In particular, even though it is trained using $\Omega_1, \ldots, \Omega_L$, when used as a reconstruction map one only has access to data from one sampling pattern of size m, and not all L sampling patterns used in the training.

This type of training does nothing to obviate instabilities, and in fact, may make them more likely. Training on more data, as in (3.13), is likely to improve the quality of the reconstruction map Ψ , making it easier to achieve the conditions

$$d_1(\Psi(P_{\Omega}Fx), x) < \eta, \quad d_1(\Psi(P_{\Omega}Fx'), x') < \eta, \quad d_1(x, x') \gg \eta.$$

Yet, the amount and variety of the training data is completely unrelated to the null space of $P_{\Omega}F$, thus it does nothing to mitigate against the condition

$$d_2(P_{\Omega}Fx, P_{\Omega}Fx') \leqslant \eta.$$

Note that the network used in Figs. 2, 4 and 8 is trained with random sampling as discussed above, yet it is highly unstable.

Adversarial training/augmenting the training set

In image classification a common strategy to enhance robustness to perturbations is to perform adversarial training [23, 33, 46]. One may view this as a way of increasing the size of the training set. Other strategies to do this include data augmentation.

However, there is no reason why an increase in the amount of training data will mitigate against instabilities in inverse problems. Suppose A is the matrix of the inverse problem (1.1). As long as the class $\mathcal{N}\mathcal{N}$ of neural networks is rich enough to ensure a small training error, then the trained network Ψ will satisfy

$$d_1(\Psi(Ax), x) < \eta, \quad d_1(\Psi(Ax'), x') < \eta,$$
 (3.14)

for some small $\eta > 0$ and all (Ax, x), (Ax', x') in the training set. However, the size of the training set is irrelevant. If

$$d_2(Ax, Ax') \leqslant \eta, \quad d_1(x, x') \gg \eta, \tag{3.15}$$

for any two such pairs (Ax, x), (Ax', x'), then Theorem 3.1 applies. Moreover, (3.15) may even be encouraged by more training data, since there are simply more pairs of (Ax, x), (Ax', x') available that satisfy (3.14).

Adding random noise

In [7, p. 138], the prospect of adding additional random noise to the measurements has been raised as a potential way to combat instabilities. This is a tempting idea, and it would have succeeded had the collection of 'bad' perturbation belonged to a set of measure zero. However, as Theorem 3.1 reveals this is not the case: the 'bad' perturbations have balls around them containing further 'bad' perturbations. Recall also the discussion in §3.2 on probabilistic aspects of instabilities. This phenomenon is illustrated in Fig. 6, where small random noise is added to the perturbation without having any mitigating effect.

4 Sparse regularization: stability through kernel awareness

Sparse regularization techniques form the core of many start-of-the-art reconstruction algorithms for inverse problems in imaging. As we explain in this section, the underlying theory for sparse regularization, *compressed sensing*, has at its heart a notion of kernel awareness. In contrast to the issues observed in the previous section, this property is sufficient to guarantee accuracy and stability of the reconstruction.

4.1 Sparse regularization

In sparse regularization, one seeks to exploit the inherent sparsity of natural images under certain sparsifying transforms, for instance, a wavelet or discrete gradient transform. If $x \in \mathbb{C}^N$ is an image and $H \in \mathbb{C}^{N \times N}$ is such a transform, then one seeks an approximate solution \widehat{x} of the linear system Ax = y for which Hx is approximately sparse. This can be achieved, for instance, by solving the quadratically-constrained basis pursuit problem

$$\widehat{x} \in \underset{z \in \mathbb{C}^N}{\operatorname{argmin}} \|Hx\|_{\ell^1} \text{ subject to } \|Az - y\|_2 \leqslant \eta.$$

Here y are the measurements of x and $\eta \geqslant 0$ is a fitting parameter, typically chosen according to the noise level of the measurements. One may also consider alternatives such as the LASSO problem

$$\widehat{x} \in \underset{z \in \mathbb{C}^N}{\operatorname{argmin}} \|Az - y\|_{\ell^2}^2 + \lambda \|Hx\|_{\ell^1},$$

where $\lambda \geqslant 0$ is a regularization parameter. Note that when H is the discrete gradient operator, $\|Hx\|_{\ell^1} = \|x\|_{TV}$ is the total variation semi-norm.

4.2 Stability and accuracy through kernel awareness

Compressed sensing theory, see [21] for an in-depth introduction, guarantees performance of sparse regularization. At its heart lies the following particular notion of kernel awareness:

Definition 4.1 (Robust Null Space Property). A matrix $A \in \mathbb{C}^{m \times N}$ satisfies the robust Null Space Property (rNSP) of order $1 \leq s \leq N$ with constants $0 < \rho < 1$ and $\gamma > 0$ if

$$||P_{\Omega}x||_{\ell^2} \leqslant \frac{\rho}{\sqrt{s}} ||P_{\Omega}^{\perp}x||_{\ell^1} + \gamma ||Ax||_{\ell^2},$$

for all $x \in \mathbb{C}^N$ and $\Omega \subseteq \{1, \dots, N\}$ with $|\Omega| \leqslant s$.

Recall that a vector $x \in \mathbb{C}^N$ is s-sparse if it has at most s nonzero components. A consequence of the rNSP is that no s-sparse vector can be close to $\mathcal{N}(A)$. Indeed, the rNSP implies that

$$||x||_{\ell^2} \leqslant \gamma ||Ax||_{\ell^2}$$
, for all s-sparse vectors x.

An implication, as we see next, is that sparse regularization achieves stable recovery when the domain \mathcal{M}_1 is taken to be the set of vectors x for which Hx is s-sparse. Moreover, we also obtain accurate recovery for vectors x that are close to \mathcal{M}_1 :

Theorem 4.2 (rNSP implies stable and robust recovery). Suppose that $H \in \mathbb{C}^{N \times N}$ is unitary and $A \in \mathbb{C}^{m \times N}$ is such that AH^* has the rNSP of order s with constants $0 < \rho < 1$ and $\gamma > 0$. Let $x \in \mathbb{C}^N$ and $y = Ax + e \in \mathbb{C}^m$, where $||e||_{\ell^2} \leq \eta$ for some $\eta \geq 0$. If

$$\widehat{x} \in \underset{\leftarrow}{\operatorname{argmin}} \|Hz\|_{\ell^{1}} \text{ subject to } \|Az - y\|_{2} \leqslant \eta, \tag{4.1}$$

then there are constants C_1, C_2 depend on ρ and γ only such that

$$\|\widehat{x} - x\|_{\ell^2} \leqslant C_1 \frac{\sigma_s(Hx)_{\ell^1}}{\sqrt{s}} + C_2 \eta,$$

where $\sigma_s(z)_{\ell^1} = \inf\{\|z - v\|_{\ell^1} : v \text{ is } s\text{-sparse}\}.$

See, for example, [21]. This theorem ensures stability – any perturbation in y of magnitude η yields an error in the recovery of x of magnitude a constant times η – as well as accuracy in the recovery of x. Indeed, the domain $\mathcal{M}_1 = \{z : Hz \text{ is } s\text{-sparse}\}$ and $\sigma_s(Hx)_{\ell^1}$ is a measure of how close x is to \mathcal{M}_1 .²

Note that this theorem applies when H is unitary, as is the case for orthogonal wavelet decompositions. Similar results can be shown for certain non-unitary cases, in particular for TV minimization (when H is the discrete gradient operator) [36].

While Theorem 4.2 assures stability, it does so only when the rNSP holds. If this condition fails to hold, one could still encounter instabilities. However, this is exactly the point: without any awareness of the kernel of A a recovery method will likely be susceptible to instabilities.

5 What do we try to learn? The optimal map

We now progress to the concept of optimal maps for inverse problems. As seen above, sparse regularization aims to recover images that belong to the domain \mathcal{M}_1 of images that are sparse under a given transform. The concept of an optimal map is motivated by the following question: given an inverse problem and a domain \mathcal{M}_1 , what is the best reconstruction possible? Indeed, this is the mapping we are aiming to construct, and the motivation for using DL is that it may provide a better reconstruction than existing sparse regularization methods where \mathcal{M}_1 is specified a priori.

5.1 Optimal maps and approximate optimal maps

So as to to achieve full generality, we need to consider multivalued maps. This is in order to include existing approaches, such as sparse regularization, that rely on minimisers of convex optimization problems which need not be unique. Recall that a multivalued mapping is traditionally noted with double arrows as

$$\varphi: \mathcal{M}_2 \rightrightarrows \mathcal{M}_1,$$

where in our cases \mathcal{M}_2 and \mathcal{M}_1 are metric spaces. We assume that the set $\varphi(x)$ is bounded for all $x \in \mathcal{M}_2$. Thus, to measure distance between $X, Z \subset \mathcal{M}_1$ we use the Hausdorff metric on the set of bounded subsets of \mathcal{M}_1 ,

$$d_{\mathcal{M}_1}^H(Z,X) = \max \big\{ \sup_{x \in X} \inf_{z \in Z} d_{\mathcal{M}_1}(z,x), \sup_{z \in Z} \inf_{x \in X} d_{\mathcal{M}_1}(z,x) \big\},$$

 $^{^2}$ A word on terminology. It is common in compressed sensing literature to refer to 'stability' and 'robustness', with the former pertaining to the effect of perturbations of x (what we refer to as 'accuracy' in this paper) and the latter refers to perturbations of y (what we refer to as 'stability').

With slight abuse of notation we will denote a singleton $\{x\} \subset \mathcal{M}_1$ by x, and we notice that $d_{\mathcal{M}_1}^H(Z,x)$ is an upper bound on the largest possible distance between x and any point in Z.

Definition 5.1 (Optimal map). Let $A: \mathbb{C}^N \to \mathbb{C}^m$ be linear, $\mathcal{M}_1 \subseteq \mathbb{C}^N$ be a metric space and $\mathcal{M}_2 = A(\mathcal{M}_1)$. Define the optimality constant

$$c_{\text{opt}}(A, \mathcal{M}_1) = \inf_{\varphi: \mathcal{M}_2 \rightrightarrows \mathcal{M}_1} \sup_{x \in \mathcal{M}_1} d_1^H(\varphi(Ax), x). \tag{5.1}$$

Since the infimum may not be attained we define an approximate optimal map as follows. We say that $\varphi_{\varepsilon}: \mathcal{M}_2 \rightrightarrows \mathcal{M}_1, \varepsilon \in (0,1]$ is a family of approximate optimal maps for (A, \mathcal{M}_1) if

$$\sup_{x \in \mathcal{M}_1} d_1^H(\varphi_{\varepsilon}(Ax), x) \leqslant c_{\text{opt}}(A, \mathcal{M}_1) + \varepsilon, \tag{5.2}$$

and that φ_0 is an optimal map if φ_0 satisfies (5.2) with $\varepsilon = 0$.

With this definition in hand, the question is now whether or not DL provides a reconstruction mapping that is either optimal or approximately optimal.

5.2 A paradox: in inverse problems we do not know the optimal map

In order to put DL for inverse problems in perspective let us recall one of the applications where DL demonstrated unprecedented success: namely, the classification problem. For simplicity we consider the simple binary classification problem

$$f: \mathcal{M} \subset \mathbb{R}^d \to \{0, 1\},$$

where \mathcal{M} is some non-empty subset. The objective of DL is to approximate f with a neural network. This is done as follows. First, one samples the graph

$$\{(x, f(x)) : x \in \mathcal{M}\},\$$

at finitely-many points $\mathcal{T} = \{(x_i, f(x_i))\}_{i=1}^K$, yielding the training set. Second, this data is fed into an optimization problem which returns a neural network $\Psi \approx f$. The crucial point is that we have access (through the training set) to the function f, which is the map we are seeking to approximate and, by definition, optimal.

The situation in inverse problems is starkly different: the optimal reconstruction map, which we are trying to learn, is not accessed through the training data. In DL for inverse problems we are given the training set

$$\mathcal{T} = \{ (y^j, x^j) : y^j = Ax^j, \ j = 1, \dots, K \}.$$
 (5.3)

and we construct an approximation of the finitely-sample 'graph' $\{(y^j, x^j)\}_{j=1}^K$. However, as we show below, there may not be an optimal map (nor a collection of approximate optimal maps) to which this graph corresponds. In other words, any mapping f which reconstructs the training set well, i.e.

$$f(y^j) \approx x^j, \qquad j = 1, \dots, n,$$

may not be an optimal map on \mathcal{M}_2 .

When compared with the classification problem, we now have the following paradoxical scenario:

- (i) In classification one knows and samples the optimal map f. Indeed, training data $\{(x^j, f(x^j))\}_{j=1}^K$ is actually a subset of the graph of f.
- (ii) In inverse problems an optimal map f (if it even exists) is unknown. Moreover, the training set $\{(y^j, x^j)\}_{j=1}^K$ is typically not a subset of the graph of f. Hence there is no reason why training should yield an optimal map.

This is formalized in Theorem 5.2, and its conclusion can be summarized as follows It is tempting to think that one should aim to approximate a function $f:Ax\mapsto x$ which performs well on the training set. After all, being able to reconstruct x from Ax should be the ultimate goal. However, as this theorem reveals, any attempt at approximating such a function may not yield an optimal map. The problem lies in the fact that we do not know the optimal map; in particular, we do not sample its graph. Hence the training process, which relies on data that may not represent a subset of the graph of the desired reconstruction mapping, may not bring us closer to the goal of reaching the optimal reconstruction.

Why does this paradoxical scenario arise? The answer lies with the noninvertibility of A. If $A \in \mathbb{C}^{N \times N}$ is invertible, then the optimal map is $f(y) = A^{-1}y$, and the training data (5.3) is indeed a subset of the graph of f. Unfortunately, this scenario is not relevant to most modern applications.

Theorem 5.2. Let the metrics $d_{\mathcal{M}_1}$ and $d_{\mathcal{M}_2}$ be on \mathbb{C}^N and \mathbb{C}^m be induced by norms. Let $A \in \mathbb{C}^{m \times N}$ with rank $(A) \geq 1$, where m < N, $K \in \{2, \ldots, \infty\}$, $\delta \leq 1/5$ and $\mathcal{B} \subset \mathbb{C}^N$ be the closed unit ball (with respect to $d_{\mathcal{M}_1}$). Then the following holds:

(i) (Training may not yield optimal maps). There exist uncountably many $\mathcal{M}_1 \subset \mathcal{B}$, such that for each \mathcal{M}_1 there exist uncountably many sets $\mathcal{T} \subset \mathcal{M}_2 \times \mathcal{M}_1$ with $|\mathcal{T}| = K$, where $\mathcal{M}_2 = A(\mathcal{M}_1)$, satisfying the following. Any map $\Psi : \mathcal{M}_2 \to \mathcal{M}_1$ (potentially multivalued $\Psi : \mathcal{M}_2 \rightrightarrows \mathcal{M}_1$) satisfying

$$d_{\mathcal{M}_1}^H(\Psi(y), x) \leqslant \delta, \qquad \forall (y, x) \in \mathcal{T},$$
 (5.4)

is not an optimal map. Moreover, the collection of such mappings does not contain a family of approximate optimal maps. If K is finite, one can choose $|\mathcal{M}_1| = K + 1$.

(ii) (The map sought by training may not exist). There exist uncountably many domains $\mathcal{M}_1 \subset \mathcal{B}$ with $|\mathcal{M}_1| = K$ such that, with $\mathcal{M}_2 = A(\mathcal{M}_1)$, there does not exist a map $\Psi : \mathcal{M}_2 \to \mathcal{M}_1$ (nor a multivalued map $\Psi : \mathcal{M}_2 \rightrightarrows \mathcal{M}_1$) for which

$$d_{\mathcal{M}_1}^H(\Psi(y), x) \leqslant \delta, \qquad \forall (y, x) \in \mathcal{M}_2 \times \mathcal{M}_1.$$

We stress that Theorem 5.2 is not a statement about overfitting. Overfitting occurs when a network performs well on the training set, but poorly on the test set. This phenomenon is caused by the fact that the architecture of the network are fixed, and hence its ability to fit data is limited (it can fit the training set, but not the test set). Indeed, it is a classical result in approximation theory that any set of data points (e.g. the union of the training and test sets) can be interpolated by a neural network of sufficient size [35]. Similarly, any continuous function can be approximated by a large enough network. So even if the trained network would suffer from overfitting, and hence lack performance on the test set, there will exist another neural network that interpolates all data points in the training set as well as the test set. What Theorem 5.2 describes is a phenomenon that happens for all mappings. Thus, there is no restriction in the network architecture used. Moreover, according to Theorem 5.2, the size of \mathcal{M}_1 can be K+1 and therefore, finite, hence the phenomenon arises regardless of the interpolation power guaranteed by classical results.

More directly, one could simply let \mathcal{T} contain both the training sets and test sets in (i). Theorem 5.2 then says that one can have excellent performance on both the training and test sets but still be a suboptimal mapping.

Remark 5.3. The condition that $\delta \leq 1/5$ is related to the assumption that \mathcal{M}_1 is a subset of the unit ball \mathcal{B} . The theorem holds for arbitrary $\delta > 0$, provided this ball is suitably enlarged.

This phenomenon is illustrated in Fig. 7. As shown therein, the trained network performs well on the left image, which lies in the training set, but performs much worse on the right image, which is close to the training set, but not in it. Conversely, sparse regularization performs well on both images.

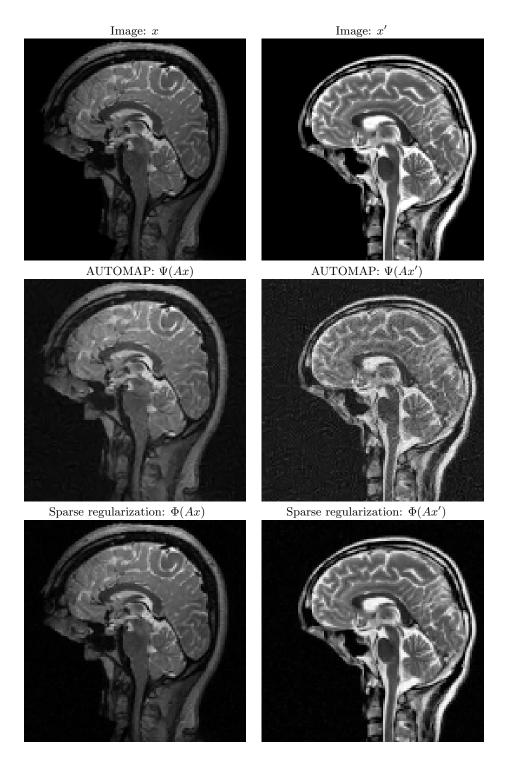


Figure 7: **(Non-optimal maps).** The AUTOMAP network [50] does not yield an optimal map. Top: Images x (left) and x' (right). The image x is from the the MGH-USC HCP public dataset [18] which was used for training and testing AUTOMAP, while the image x' is from General Electric Healthcare (and belongs to a completely different dataset). Middle: Reconstructions of x (left) and x' (right) from measurements Ax and Ax' using the AUTOMAP network $\Psi: \mathbb{C}^m \to \mathbb{C}^N$. Bottom: Reconstructions of x (left) and x' (right) from measurements Ax and Ax' using a sparse regularization decoder $\Phi: \mathbb{C}^m \to \mathbb{C}^N$. As is evident, the AUTOMAP network struggles to recover the new image x' well.

6 Deep learning or sparse regularization?

Since sparse regularization has over the last decade frequently become the method of choice for imaging, it is natural to examine how it compares with DL techniques for different choices of \mathcal{M}_1 . In this section, we give some theoretical insights into their comparative performance. In particular, we will demonstrate situations where sparse regularization is outperformed by DL, but only at the cost of instability, and where sparse regularization and DL may both be optimal, but the latter may be destabilized by adding more training data

Throughout this section the norm used, unless otherwise stated, will be the ℓ^2 -norm. Also, unless otherwise stated, the metric on any metric space will be that induced by the ℓ^2 -norm. In particular, we write simply d^H for the corresponding Hausdorff metric. However, the choice of norm is not particularly important as different choices would only change the constants in main results. For convenience, we now also define the ε -Lipschitz constant of a mapping $\Phi: \mathbb{C}^m \rightrightarrows \mathbb{C}^N$ over a set $\mathcal{M}_2 \subseteq \mathbb{C}^m$:

$$L^{\varepsilon}(\Phi, \mathcal{M}_2) = \sup_{y \in \mathcal{M}_2, 0 < \|e\| \leqslant \varepsilon} \frac{d^H(\Phi(y+e), \Phi(y))}{\|e\|}.$$

6.1 Compressed sensing for inverse problems in imaging

The classical theory of compressed sensing, where \mathcal{M}_1 is the set of s-sparse vectors (in some transform), is typically insufficient for image reconstruction problems [4, 37, 38]. For such problems, one needs a more refined model that incorporates the *local* properties of the sparsifying transform and the sampling operator (e.g. the Fourier transform). This was developed in a series of works [1, 2, 3, 9, 29]

We now describe this setup. Throughout, to keep the level of technicality as low as possible, we consider one-dimensional, discrete images of size $N=2^{r-1},\,r\geqslant 2$, and let $H\in\mathbb{C}^{N\times N}$ be the discrete Haar wavelet transform. Neither assumption is fundamental: both higher dimensions and other wavelets can also be considered $[1,\,2]$.

Images, when sparsified with a wavelet sparsifying transform, possess coefficients that are not only approximately sparse, but have a specific sparsity structure within the wavelet scales. This motivates the *sparsity in levels* model:

Definition 6.1 (Sparsity in levels). Let $\mathbf{M} = (M_1, \dots, M_r) \in \mathbb{N}^r$, $M_0 = 0$, with $1 \leq M_1 < \dots < M_r = N$ and let $\mathbf{s} = (s_1, \dots, s_r) \in \mathbb{N}^r$ with $s_l \leq M_l - M_{l-1}$, for $l = 1, \dots, r$. A vector $x \in \mathbb{C}^N$ is (\mathbf{s}, \mathbf{M}) -sparse in levels if

$$|\text{Supp}(x) \cap \{M_{l-1} + 1, \dots, M_l\}| \le s_l \text{ for } l = 1, \dots, r.$$

The set of all (s, \mathbf{M}) -sparse vectors is denoted by $\Sigma_{s, \mathbf{M}}$.

While the above definition is fully general, in what follows the levels \mathbf{M} will be chosen to delineate the Haar wavelet scales. Specifically, $M_l = 2^{l-1}$ for $l = 1, \ldots, r$, so that the first level contains the scaling function coefficient, and the l^{th} level, $l \geq 2$, contains the wavelet coefficients at scale l-2.

We now turn our attention to the measurements. There are two main types of sampling that dominate compressed sensing in image reconstruction. First, sampling with the Fourier transform, as is the case in, for instance, MRI. Second, binary sampling with the Walsh–Hadamard transform, as is the case in numerous optical imaging modalities. See [1] and references therein. The success of sparse regularization in imaging owes itself to the careful design, via compressed sensing theory, of sampling patterns that are tailored to the sparsity in levels structure of the images to be recovered (see [4, 38] and references therein). Note that a sampling pattern is a choice of m Fourier or Walsh frequencies selected amongst all N possible frequencies. Specifically, we consider the following types of sampling patterns:

Definition 6.2 (Multilevel random sampling). Let $\mathbf{N} = (N_1, \dots, N_r) \in \mathbb{N}^r$, where $1 \leq N_1 < \dots < N_r = N$ and $\mathbf{m} = (m_1, \dots, m_r) \in \mathbb{N}^r$ with $m_k \leq N_k - N_{k-1}$ for $k = 1, \dots, r$, and

 $N_0 = 0$. For each k = 1, ..., r, let $\Omega_k = \{N_{k-1} + 1, ..., N_k\}$ if $m_k = N_k - N_{k-1}$ and if not, let $t_{k,1}, ..., t_{k,m_k}$ be chosen uniformly and independently from the set $\{N_{k-1} + 1, ..., N_k\}$, and set $\Omega_k = \{t_{k,1}, ..., t_{k,m_k}\}$. We refer to $\Omega = \Omega_{\boldsymbol{m},\mathbf{N}} = \Omega_1 \cup \cdots \cup \Omega_r$ as an $(\boldsymbol{m},\mathbf{N})$ -multilevel random sampling scheme.

The values **N** are referred to as sampling levels. The total number of measurements is $m=m_1+\ldots+m_r$. Let $U\in\mathbb{C}^{N\times N}$ be a fixed, unitary matrix. Then a multilevel random sampling scheme gives rise to a measurement matrix $A\in\mathbb{C}^{m\times N}$ of the form

$$A = P_{\Omega}DU, \qquad D = \operatorname{diag}(d_1, \dots, d_N), \tag{6.1}$$

where

$$d_i = \begin{cases} \sqrt{\frac{N_k - N_{k-1}}{m_k}} & \text{if } m_k \neq 0 \\ 1 & \text{if } m_k = 0 \end{cases}, \ N_{k-1} < i \leqslant N_k, \ k = 1, \dots, r,$$

(the diagonal matrix D is a scaling matrix to ensure that $\mathbb{E}(A^*A) = I$). While general, in what follows we take U to be either a discrete Fourier or Walsh–Hadamard transform, and choose the sampling levels as $N_k = 2^{k-1}$, $k = 1, \ldots, r$. The first level then corresponds to the lowest Fourier/Walsh frequency, and the l^{th} level, $l \geq 2$, contains frequencies in a dyadic band of size roughly 2^l .

Choosing sparsity in levels as our model, we require a slightly different regularization term. Following [1], we replace the ℓ^1 -norm by a weighted ℓ^1 -norm, where the weights are chosen to be constant on the sparsity levels. We define

$$||x||_{\ell^1_{\boldsymbol{\omega}}} = \sum_{l=1}^r \omega_l ||P_{M_l}^{M_{l-1}} x||_{\ell^1}, \quad x \in \mathbb{C}^N,$$

where $\boldsymbol{\omega} = (\omega_1, \dots, \omega_r)$ is a set of positive weights, and $P_{M_l}^{M_{l-1}} = P_{\{M_{l-1}+1,\dots,M_l\}}$ is the projection onto the l^{th} sparsity level.

Let $\mathcal{M}_1 = H^*(\Sigma_{\mathbf{s},\mathbf{M}})$ be the set of images whose Haar wavelet coefficients are (\mathbf{s},\mathbf{M}) -sparse, and set $\mathcal{M}_2 = A(\mathcal{M}_1)$. We now define the compressed sensing mapping $\Phi : \mathbb{C}^m \rightrightarrows \mathbb{C}^N$ via the following weighted ℓ^1 minimization problem

$$\Phi(y) = \underset{z \in \mathbb{C}^N}{\operatorname{argmin}} \|Hz\|_{\ell^1_{\omega}} \text{ subject to } \|Az - y\|_{\ell^2} \leqslant \eta, \tag{6.2}$$

where $\eta = \inf\{\|y - z\|_{\ell^2} : z \in \mathcal{M}_2\}$ is the ℓ^2 -norm distance from y to \mathcal{M}_2 , and $\boldsymbol{\omega} = (\sqrt{s/s_1}, \dots, \sqrt{s/s_r})$.

6.2 Compressed sensing is sometimes optimal, but may be outperformed by deep learning at the cost of instability

In this section, we consider binary sampling via the Walsh–Hadamard transform $U \in \mathbb{C}^{N \times N}$, where, as above, $N = 2^{r-1}$. The sparsity and sampling levels **M** and **N** are also as defined above. We also now use the notation $A \lesssim B$ to mean there exists a numerical constant c > 0 such that $A \leqslant cB$.

Theorem 6.3. Let $U \in \mathbb{R}^{N \times N}$ be the sequency-ordered Walsh-Hadamard transform, where $N = 2^{r-1}$, $r \geq 2$, and $0 < \nu < 1$. There exist at least $2^{(r-1)(r-2)/2}$ different (s, \mathbf{M}) sparsity patterns such that for each such pattern and each $p \in (2, \infty)$, there exist uncountably many domains $\mathcal{M}_1, \widetilde{\mathcal{M}}_1 \subset \mathbb{C}^N$ with $|\mathcal{M}_1| = |\widetilde{\mathcal{M}}_1| = \infty$ for which the following holds. If $\Omega = \Omega_{\mathbf{m}, \mathbf{N}}$ is an (\mathbf{m}, \mathbf{N}) -multi-level random sampling scheme, where

$$m_k = \left[Cs_k \left(\log^3(N) \log(2m) \log^2(2s) + \log(N) \log(\nu^{-1}) \right) \right], \quad k = 1, \dots, r,$$
 (6.3)

C > 0 is a universal constant and $m = m_1 + \ldots + m_r$, A is as in (6.1), $\mathcal{M}_2 = A(\mathcal{M}_1)$ and $\widetilde{\mathcal{M}}_2 = A(\widetilde{\mathcal{M}}_1)$, then the following statements holds:

- (i) (CS is stable and optimal on \mathcal{M}_1). With probability greater than 1ν , the mapping Φ defined in (6.2) is an optimal mapping (see Definition 5.1) and the optimality constant satisfies $c_{\mathrm{opt}}(A,\mathcal{M}_1)=0$. Moreover, the ε -Lipschitz constant $L^{\varepsilon}(\Phi,\mathcal{M}_2)\lesssim r^{1/4}$ for all $\varepsilon>0$. Specifically, $L^{\varepsilon}(\Phi,\mathcal{M}_2)\leqslant 2\sqrt{2}+(1+r^{1/4})8\sqrt{2}/(2-\sqrt{3})$.
- (ii) (DL outperforms CS on $\widetilde{\mathcal{M}}_1$ at the cost of instability). For any $\delta \in (0, 1/p]$ and any mapping $\Psi : \widetilde{\mathcal{M}}_2 \rightrightarrows \widetilde{\mathcal{M}}_1$ satisfying

$$d^{H}(\Psi(Ax), x) \leqslant \delta, \qquad \forall x \in \widetilde{\mathcal{M}}_{1},$$
 (6.4)

it holds that

$$\sup_{x \in \widetilde{\mathcal{M}}_1} d^H(\Psi(Ax), x) \leqslant \frac{1}{p} \sup_{x \in \widetilde{\mathcal{M}}_1} d^H(\Phi(Ax), x), \tag{6.5}$$

i.e. Ψ is at least p times closer to an optimal map than Φ . However,

$$L^{\varepsilon}(\Phi, \widetilde{\mathcal{M}}_2) \leqslant 2\sqrt{2} + (1 + r^{1/4})8\sqrt{2}/(2 - \sqrt{3}) \quad \forall \varepsilon > 0, \tag{6.6}$$

with probability greater than $1 - \nu$, whereas

$$L^{\varepsilon}(\Psi, \widetilde{\mathcal{M}}_2) \geqslant \frac{1}{\varepsilon}, \quad \varepsilon \geqslant \frac{1}{p}.$$
 (6.7)

In other words, the instability of Ψ increases with its performance.

Theorem 6.3 demonstrates how compressed sensing may be optimal on one domain, but, on another domain DL will provably be better. Moreover, the improvement can be arbitrarily good, as shown in (6.5). However, this comes at the cost of instability. Indeed, as (6.5) and (6.6) demonstrate, the better the performance gets the more unstable the NN mapping becomes. Conversely, the Lipschitz constant of the CS mapping is at most $\mathcal{O}\left(r^{1/4}\right) = \mathcal{O}\left((\log(N))^{1/4}\right)$, so for all intents and purposes, bounded. Note that the constant factors in the bound on the Lipschitz constant are not in any way special. They are related to the constant C in (6.3). Fig. 8 gives an illustration of this phenomenon.

6.3 Compressed sensing and deep learning may both be optimal, but two wrong training elements may destabilize deep learning

In this section, we consider the same setup as the previous section, except now $U \in \mathbb{C}^{N \times N}$, where $N = 2^{r-1}$, is taken as the discrete Fourier transform. The results in this section hold for classes of NNs that are rich enough so as to interpolate certain collections of finitely-many points. Specifically:

Definition 6.4 (*M*-interpolatory networks). Given $M \in \mathbb{N}$ we say that a class \mathcal{NN} of neural networks $\Phi : \mathbb{C}^m \to \mathbb{C}^N$ is M-interpolatory if for each collection $\{x_j\}_{j=1}^M \subset \mathbb{C}^N$, $\{y_j\}_{j=1}^M \subset \mathbb{C}^m$ there exists a neural network $\Psi \in \mathcal{NN}$ such that $\Psi(y_j) = x_j, \forall j = 1, \ldots, M$.

Note that the interpolation properties of one-layer neural networks are well understood [35]. Specifically, when N=1 a network with one hidden layer is M-interpolary, provided the layer has at least M nodes (in the case of real inputs and outputs) or 2M nodes (in the complex case) and the activation function is not a polynomial. The same result trivially holds for deeper networks, as long as the number of nodes in any of the hidden layers exceeds M, although this bound is likely not tight.

Having established the concept of M-interpolatory networks we can now state the main theorem of this section.

Theorem 6.5. Let $U \in \mathbb{C}^{N \times N}$ be the discrete Fourier transform, where $N = 2^{r-1}$, $r \geqslant 2$, and $0 < \nu < 1$. Let $K, M \in \mathbb{N}$ with $M \geqslant K$. For any (s, M) sparsity pattern, $s \neq 0$, there exists

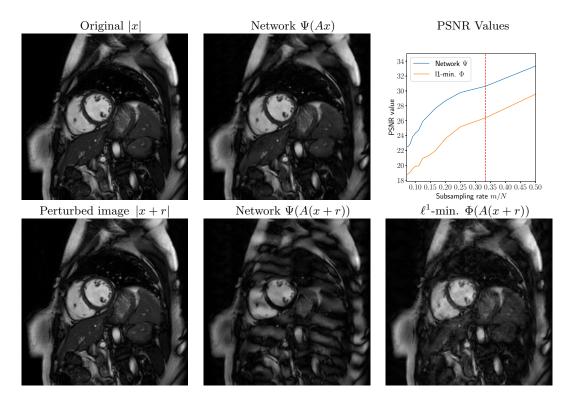


Figure 8: (Instability through overperformance.) In all experiments $A \in \mathbb{C}^{m \times N}$ is a subsampled Fourier transform. Left: (top) Original image |x| and (bottom) perturbed image |x+r|. Middle: Reconstruction of the original and perturbed images, x and x+r from measurements Ax (top) and A(x+r) (bottom) using the deep MRI network $\Psi \colon \mathbb{C}^m \to \mathbb{C}^N$ with m/N = 0.33. Right: (bottom) reconstruction from A(x+r) with m/N = 0.33 using the reconstruction map $\Phi \colon \mathbb{C}^m \to \mathbb{C}^N$ obtained using standard ℓ^1 -minimization with wavelets. In the top right image we plot PSNR-values against subsampling rate for the two reconstruction maps Ψ and Φ . We choose different sampling rates m/N in the interval $\left[\frac{1}{14},\frac{1}{2}\right]$ and preform reconstruction of the image x from measurements Ax. The red dashed line indicates at which sampling rate the network is trained.

uncountably many sets $\{x_1, \ldots, x_K\}$, of size K, and domains $\mathcal{M}_1 \supset \{x_1, \ldots, x_K\}$ of size M for which the following holds. If $\Omega = \Omega_{\mathbf{m}, \mathbf{N}}$ is a multi-level random sampling scheme satisfying

$$m_k = \left[C \left(s_k + \sum_{l=1}^{k-1} s_l 2^{-(k-l)} + \sum_{l=k+1}^r s_l 2^{-3(l-k)} \right) \cdot L \right], \quad k = 1, \dots, r,$$
 (6.8)

where C > 0 is a universal constant,

$$L = \log^3(N) \cdot \log(2m) \cdot \log^2(2s) + \log(N) \cdot \log(\nu^{-1}),$$

and $m = m_1 + \ldots + m_r < N$, A is as in (6.1) and $\mathcal{M}_2 = A(\mathcal{M}_1)$, then statements (i) and (ii) below hold.

- (i) (CS is stable and optimal on \mathcal{M}_1). With probability greater than 1ν , the mapping Φ defined in (6.2) is an optimal map where the optimality constant $c_{\text{opt}}(A, \mathcal{M}_1) = 0$, and the ε -Lipschitz constant $L^{\varepsilon}(\Phi, \mathcal{M}_2) \lesssim r^{1/4}$ for all $\varepsilon > 0$. Specifically, $L^{\varepsilon}(\Phi, \mathcal{M}_2) \leqslant 2\sqrt{2} + (1 + r^{1/4})8\sqrt{2}/(2 \sqrt{3})$.
- (ii) (DL is also optimal on \mathcal{M}_1 , but additional training data may destabilize). Let

 $\mathcal{N}\mathcal{N}$ be $\max\{M, K+1\}$ -interpolatory. Then there exists a neural network Ψ such that

$$\Psi \in \operatorname*{argmin}_{\widetilde{\Psi} \in \mathcal{NN}} \frac{1}{K} \sum_{i=1}^{K} \|x_i - \widetilde{\Psi}(Ax_i)\|_{\ell^2}^2,$$

and Ψ is an optimal map on \mathcal{M}_1 . However, for any $\gamma \in (0,1)$ there are uncountably many sets $\{x_{K+1}, x_{K+2}\}$, where $\|x_{K+1} - x_j\|$, $\|x_{K+2} - x_i\| \leq \gamma$ for some $i, j \leq K$ such that any neural network

$$\widehat{\Psi} \in \operatorname*{argmin}_{\widetilde{\Psi} \in \mathcal{NN}} \frac{1}{K+2} \sum_{i=1}^{K+2} \frac{1}{2} \|x_i - \widetilde{\Psi}\left(Ax_i\right)\|_{\ell^2}^2,$$

has an ε -Lipschitz constant bounded from below by

$$L^{\varepsilon}(\widehat{\Psi}, \mathcal{M}_2) \geqslant \frac{1}{\varepsilon}, \quad \varepsilon \geqslant \gamma.$$

Theorem 6.5 illustrates how fragile the training process for NNs can be. Indeed, by adding two extra elements in the training set the ε -Lipschitz constant can become arbitrary large. Moreover, these new elements can be arbitrary close to any of the elements in the original training set.

7 Stability versus performance: Setting the regularization parameter is challenging

Creating a stable neural network for inverse problems is not hard. Indeed, one may simply consider the zero network, a highly stable, yet not entirely useful network. The challenge is combining stability with performance. Note that Theorem 3.1 reveals that it is the current training process that causes the instabilities. Indeed, solving, for example,

$$\Psi \in \underset{\widetilde{\Psi} \in \mathcal{NN}}{\operatorname{argmin}} \frac{1}{K} \sum_{j=1}^{K} \operatorname{Cost}(\widetilde{\Psi}, Ax^{j}, x^{j}), \tag{7.1}$$

where Cost is an appropriate cost function, for example

$$\operatorname{Cost}(\widetilde{\Psi}, Ax_i, x_i) = \|x_i - \widetilde{\Psi}(Ax_i)\|_{\ell^2}^2$$

produces a network with small training error

$$\|\Psi(y) - x\| \le \delta, \qquad \forall (y, x) \in \mathcal{T},$$
 (7.2)

for some $\delta > 0$ and some norm $\|\cdot\|$, where \mathcal{T} denotes the training set. However, as Theorem 5.2 reveals, it is not just instabilities that are caused by (7.2). Indeed, this property also creates suboptimal maps.

Although, as noted, these issues are not the same as overfitting (see §5), one may attempt to overcome them in a similar way by adding a regularization term to the objective function (see §2). Specifically, one solves

$$\Psi \in \underset{\widehat{\Psi} \in \mathcal{N} \mathcal{N}}{\operatorname{argmin}} \frac{1}{|\mathcal{T}|} \sum_{(y,x) \in \mathcal{T}} \frac{1}{2} \|x - \widehat{\Psi}(y)\|_{\ell^{2}}^{2} + \lambda J(\widehat{\Psi})$$

$$(7.3)$$

where $\lambda \in \mathbb{R}_+ := [0, \infty)$, and $J : \mathcal{NN} \to \mathbb{R}$ is a function from a class of neural networks to the reals. This raises the question: how does one set the regularisation parameter λ in (7.3)? In this section, we investigate. To do so, we need the following:

Definition 7.1 (Optimal λ). Given a pair (A, \mathcal{M}_1) , a class $\mathcal{N}\mathcal{N}$ of neural networks and a training set \mathcal{T} we say that $\lambda \in \mathbb{R}_+$ is optimal for $\{(A, \mathcal{M}_1), \mathcal{N}\mathcal{N}, \mathcal{T}\}$ if there is a minimiser of (7.3) that is an optimal map for (A, \mathcal{M}_1) .

As Theorem 7.2 reveals setting the right λ is a highly delicate, and ironically, highly unstable problem:

Theorem 7.2 (Setting λ is delicate). Let $U \in \mathbb{C}^{N \times N}$ be an invertible matrix with $N \geqslant 4$. There exist $2^N - 2N - 2$ sampling patterns $\widetilde{\Omega}$ and for each of them a sampling pattern Ω , that is not full sampling, such that if $A = P_{\Omega}U$ and $\widetilde{A} = P_{\widetilde{\Omega}}U$ we have the following. Let $K \in \{2, \ldots, \infty\}$, $\mathcal{N}\mathcal{N} = \mathcal{N}\mathcal{N}_{\mathbf{n}}$ be any class of ReLU neural networks with at least one hidden layer, where $\mathbf{n} = (m, n_1, \ldots, n_{L-1}, N)$, $m = \max\{|\Omega|, |\widetilde{\Omega}|\}$ (see Remark 7.3), $n_j \geqslant 2m$, $(n_j \geqslant 4m$, if we consider it as a real-valued network, see Remark 2.1) and $J : \mathcal{N}\mathcal{N} \to \mathbb{R}$. Then, there is a $\lambda_{\mathrm{opt}} \in \mathbb{R}_+$ and uncountably many domains \mathcal{M}_1 of size K such that for each \mathcal{M}_1 there are uncountably many training sets $\mathcal{T} = \mathcal{R} \times \mathcal{D} \subset A(\mathcal{M}_1) \times \mathcal{M}_1$ and $\widetilde{\mathcal{T}} = \widetilde{\mathcal{R}} \times \widetilde{\mathcal{D}} \subset \widetilde{A}(\mathcal{M}_1) \times \mathcal{M}_1$ of finite size and possibly with $\mathcal{D} = \widetilde{\mathcal{D}}$, such that

$$\lambda_{\text{opt}}$$
 is optimal for $\{(\widetilde{A}, \mathcal{M}_1), \mathcal{NN}, \widetilde{\mathcal{T}}\}\$ and $\{(A, \mathcal{M}_1), \mathcal{NN}, \mathcal{T}\}.$

However, there exists an uncountable collection $S \subseteq \mathbb{C}^N$ such that the following holds:

(i) (The sampling pattern $\widetilde{\Omega}$ makes λ_{opt} unstable with respect to \mathcal{S}) If $(\widetilde{A}x, x), x \in \mathcal{S}$, is either added to the training set $\widetilde{\mathcal{T}}$, or replaces a specific element in $\widetilde{\mathcal{T}}$, then either there is no element of

$$\underset{\widehat{\Psi} \in \mathcal{NN}}{\operatorname{argmin}} \frac{1}{|\widetilde{\mathcal{T}}|} \sum_{(\widetilde{y}, \widetilde{x}) \in \widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \widehat{\Psi}(\widetilde{y})\|_{\ell^{2}}^{2} + \lambda J(\widehat{\Psi}), \tag{7.4}$$

for any $\lambda \in \mathbb{R}_+$ that is an optimal map for $(\widetilde{A}, \mathcal{M}_1)$, or there is another $\widetilde{\lambda}_{\mathrm{opt}} \neq \lambda_{\mathrm{opt}}$ that is optimal for $\{(\widetilde{A}, \mathcal{M}_1), \mathcal{NN}, \widetilde{\mathcal{T}}\}$ whereas λ_{opt} is not.

(ii) (The sampling pattern Ω makes λ_{opt} stable with respect to S) Yet, given any subset $V \subset S$ such that if $\{(Ax, x) \mid x \in V\}$ is either added to T or replaces elements in T, then λ_{opt} is still optimal for $\{(A, \mathcal{M}_1), \mathcal{NN}, \mathcal{T}\}$.

Remark 7.3. Note that in the case $|\Omega| \neq |\Omega|$ and $m = \max\{|\Omega|, |\Omega|\}$ we interpret the action of $\Psi \in \mathcal{NN}_{\mathbf{n}}$ on y denoted, with slight abuse of notation, as $\Psi(y)$ as Ψ acting on $y \oplus 0$ if y has dimension less than m.

Theorem 7.2 says that given an invertible matrix U, there is an abundance of sampling patterns Ω and $\widetilde{\Omega}$, as well as training sets \mathcal{T} and $\widetilde{\mathcal{T}}$ plus a large set \mathcal{S} , such that setting the optimal parameter λ is highly unstable with respect to changes in the training set $\widetilde{\mathcal{T}}$ from elements in \mathcal{S} when considering the sampling pattern $\widetilde{\Omega}$. However, at the same time, setting the optimal parameter λ is highly stable with respect to changes in the training set \mathcal{T} from the same elements in \mathcal{S} when considering the sampling pattern Ω . Hence, changes in the training sets from the same collection \mathcal{S} can give vastly different results. The conclusion is therefore that unless one has prior information about the training data, or a potential way of learning this information, setting the λ parameter is a delicate affair. Ironically, one ends up with a potentially unstable problem in order to cure the instability issue in the original problem.

Acknowledgments

NMG acknowledges support from a UK EPSRC grant. BA acknowledges the support of the PIMS CRG "High-dimensional Data Analysis", SFU's Big Data Initiative "Next Big Question" Fund and NSERC through grant R611675. ACH acknowledges support from a Royal Society University Research Fellowship and the Leverhulme Prize 2017.

8 Proofs of the main results

We have divided the proofs into five subsections, one of which contains only preliminary results from the compressed sensing literature. All results are proved in the order they appear in the main paper.

8.1 Proofs of Theorem 3.1, Corollary 3.3 and Theorem 5.2

Proof of Theorem 3.1. For part (i), note that, by the definition of the local ε -Lipschitz constant, by assumption (3.2) and by applying the reverse triangle inequality twice, we get

$$\begin{split} L^{\varepsilon}(\Psi,y) &= \sup_{0 < d_2(z,y) \leqslant \varepsilon} \frac{d_1(\Psi(z),\Psi(y))}{d_2(z,y)} \geqslant \frac{d_1(\Psi(Ax'),\Psi(Ax))}{d_2(Ax',Ax)} \\ &\geqslant \frac{d_1(x',x) - d_1(\Psi(Ax),x) - d_1(\Psi(Ax'),x')}{d_2(Ax',Ax)}. \end{split}$$

Thus, using the fact that $\varepsilon \geqslant \eta$ and assumption (3.1) it follows that

$$L^{\varepsilon}(\Psi, y) > \frac{d_1(x', x) - 2\eta}{\varepsilon}.$$

The continuity of Ψ now implies that

$$L^{\varepsilon}(\Psi, \widetilde{y}) \geqslant \frac{1}{\varepsilon} (d_1(x, x') - 2\eta),$$

for all $\widetilde{y} \in \mathcal{B}$ where \mathcal{B} is some ball around y = Ax.

For part (ii), let z = x' - x and e = A(x' - x) and observe that, by assumption (3.2) and the fact that the metrics are translation invariant, we have $d_1(0, z) \ge d_1(x, x')$, $d_2(0, e) \le \eta$. Moreover,

$$d_1(\Psi(Ax+e), x+z) = d_1(\Psi(Ax'), x') < \eta,$$

by assumption (3.1). The result once more follows from continuity of Ψ . Part (iii) is identical.

Proof of Corollary 3.3. Let \widetilde{y} with $\|y-\widetilde{y}\|_{\ell^2} \leqslant \sigma_{\min}(A)r_1$ (this gives the closed ball \mathcal{B}_2) and write $\widetilde{x} = A^{\dagger}\widetilde{y}$, where \dagger denotes the pseudoinverse. Since $x \in \mathcal{N}(A)^{\perp}$ and A is full rank, we have $x = A^{\dagger}y$. In particular, $\|\widetilde{x} - x\|_{\ell^2} \leqslant 1/\sigma_{\min}(A)\|\widetilde{y} - y\|_{\ell^2} \leqslant r_1$. Hence $\widetilde{x} \in \mathcal{B}_1$. Set $e = A(x' - \widetilde{x})$ and observe that $\|e\|_{\ell^2} \leqslant \eta \leqslant \varepsilon$. Then, since $A\widetilde{x} = \widetilde{y}$,

$$\begin{split} L_{\varepsilon}(\Psi,\widetilde{y}) &\geqslant \frac{\|\Psi(A\widetilde{x}+e) - \Psi(A\widetilde{x})\|_{\ell^{2}}}{\|e\|_{\ell^{2}}} \\ &\geqslant \frac{\|\Psi(Ax') - \Psi(A\widetilde{x})\|_{\ell^{2}}}{\varepsilon} \\ &\geqslant \frac{\|x' - \widetilde{x}\|_{\ell^{2}} - \|x' - \Psi(Ax')\|_{\ell^{2}} - \|\widetilde{x} - \Psi(A\widetilde{x})\|_{\ell^{2}}}{\varepsilon} \\ &\geqslant \frac{\operatorname{dist}(x',\mathcal{B}_{1}) - 2\eta}{\varepsilon}. \end{split}$$

This completes the proof.

Proof of Theorem 5.2. The different norms on \mathbb{C}^N and \mathbb{C}^m are not important to distinguish in the proof. Hence we will use the notation $\|\cdot\|$ everywhere. We begin with the proof of (ii). First, let $\{x_1,\ldots,x_K\}$ be K distinct elements in $\mathcal{N}(A))^{\perp}$ such that $\|x_1\|=1/2$ and $0<\|x_j\|\leqslant 1$. Note that we can do this since, by assumption, $\operatorname{rank}(A)\geqslant 1$. Choose, by the assumption that $m< N, z_1\in\mathcal{N}(A)$ with $\|z_1\|=1/2$. Let $\mathcal{M}_1=\{x_1+z_1,x_1,\ldots,x_K\}$ and observe that

 $\mathcal{M}_1 \subset \mathcal{B}$. We argue by contradiction and suppose that there exists a (possibly multivalued) map $\Psi: \mathcal{M}_2 \rightrightarrows \mathcal{M}_1$ with

$$d_{\mathcal{M}_1}^H(\Psi(Ax), x) \leqslant \delta, \qquad \forall x \in \mathcal{M}_1.$$
 (8.1)

In particular, $d_{\mathcal{M}_1}^H(\Psi(Ax_1), x_1) \leq \delta$. However, since $z_1 \in \mathcal{N}(A)$, we have $\Psi(A(x_1+z_1)) = \Psi(Ax_1)$ and therefore $d_{\mathcal{M}_1}^H(\Psi(A(x_1+z_1)), x_1+z_1) \geq ||z_1|| - \delta \geq 1/2 - 1/5 > \delta$, which contradicts (8.1). Since K is arbitrary, we get the result. In order to get uncountably many different \mathcal{M}_1 's, as mentioned in the statement of the theorem, one can simply multiply the original choice of \mathcal{M}_1 by complex numbers of modulus 1.

To prove (i) we use the setup from the proof of (ii). Indeed, let \mathcal{M}_1 be as defined previously, and set $\mathcal{T} = \{x_1, \dots, x_K\}$. Define the map $\psi_0 : \mathcal{M}_2 \to \mathcal{M}_1$ by

$$\psi_0(y) = \begin{cases} x_1 + \frac{1}{2}z_1 & \text{if } y = Ax_1 \\ x_1 + \frac{1}{2}z_1 & \text{if } y = A(x_1 + z_1) \\ x_j & \text{otherwise} \end{cases}$$
 (8.2)

Then, by (8.2),

$$c_{\text{opt}}(A, \mathcal{M}_{1}) = \inf_{\varphi: \mathcal{M}_{2} \rightrightarrows \mathcal{M}_{1}} \sup_{x \in \mathcal{M}_{1}} d_{1}^{H}(\varphi(Ax), x) \leqslant \sup_{x \in \mathcal{M}_{1}} \|\psi_{0}(Ax) - x\|$$

$$\leqslant \sup_{j \geqslant 1} \|\psi_{0}(Ax_{j}) - x_{j}\| \vee \|\psi_{0}(A(x_{1} + z_{1})) - (x_{1} + z_{1})\| = \frac{1}{4},$$
(8.3)

where $a \vee b$ denotes the standard maximum of real numbers a, b. However, for any mapping $\Psi : \mathcal{M}_2 \rightrightarrows \mathcal{M}_1$ with

$$d_{\mathcal{M}_1}^H(\Psi(Ax_j), x_j) \leqslant \delta, \quad \forall j = 1, \dots, K,$$
(8.4)

we have that

$$\sup_{x \in \mathcal{M}_1} d_{\mathcal{M}_1}^H(\Psi(Ax), x) \geqslant d_{\mathcal{M}_1}^H(\Psi(A(x_1 + z_1)), x_1 + z_1)
= d_{\mathcal{M}_1}^H(\Psi(Ax_1), x_1 + z_1) \geqslant ||z_1|| - d_{\mathcal{M}_1}^H(\Psi(Ax_1), x_1) \geqslant \frac{1}{2} - \frac{1}{5} = \frac{3}{10} > \frac{1}{4}.$$

Thus, by (8.3), it follows that Ψ is not an optimal map. Furthermore, it is clear that no family of maps satisfying (8.4) can be approximately optimal.

8.2 Preliminary compressed sensing results

Before we prove Theorem 6.3 and 6.5 we require several results from [1]. The following is [1, Defn. 2.8]:

Definition 8.1. Let $U \in \mathbb{C}^{N \times N}$ be a unitary matrix and **N** and **M** be sampling and sparsity levels respectively. The (k, l)th local coherence of U is

$$\mu_{k,l}(U) = \max\{|U_{ij}|^2 : N_{k-1} < i \le N_k, \ M_{l-1} < j \le M_l\}.$$

The next result is [1, Cor. 3.7], simplified to the finite-dimensional case:

Theorem 8.2. Let $U \in \mathbb{C}^{N \times N}$ and $H \in \mathbb{C}^{N \times N}$ be unitary, $\Omega = \Omega_{\mathbf{m}, \mathbf{N}}$ be an (\mathbf{m}, \mathbf{N}) -multilevel random sampling scheme with r levels and $0 < \nu < 1$. Let let $\mathbf{M}, \mathbf{s} \in \mathbb{N}^r$ be sparsity levels and local sparsities respectively and $\boldsymbol{\omega} = (\sqrt{s/s_1}, \ldots, \sqrt{s/s_r})$ be weights. Let $\varepsilon, \delta \in (0, 1)$, $m = m_1 + \ldots + m_r$, $s = s_1 + \ldots + s_r$, $L = r \cdot \log(2m) \cdot \log(2N) \cdot \log^2(2s) + \log(\nu^{-1})$ and suppose that

$$m_k \gtrsim r \cdot (N_k - N_{k-1}) \cdot \left(\sum_{l=1}^r \mu_{k,l} s_l\right) \cdot L, \quad k = 1, \dots, r,$$

where $\mu_{k,l} = \mu_{k,l}(UH^*)$. Then with probability at least $1 - \nu$ the following holds. For every $x \in \mathbb{C}^N$ and $e \in \mathbb{C}^m$ with $||e||_2 \leq \eta$, any solution \hat{x} of the optimization problem

$$\min_{z \in \mathbb{C}^N} \|Hz\|_{\ell^1_{\omega}} \text{ subject to } \|Az - y\|_{\ell^2} \leqslant \eta, \tag{8.5}$$

where y = Ax + e, satisfies

$$||x - \widehat{x}||_{\ell^2} \leqslant (1 + r^{1/4}) \left(C \frac{\sigma_{s,\mathbf{M}}(Hx)_{\ell^1_{\boldsymbol{\omega}}}}{\sqrt{rs}} + D\eta \right), \tag{8.6}$$

where $\sigma_{\mathbf{s},\mathbf{M}}(z)_{\ell^1} = \inf\{\|z - v\|_{\ell^1_{\boldsymbol{\omega}}} : v \in \Sigma_{\mathbf{s},\mathbf{M}}\}, \ C = 2(2 + \sqrt{3})/(2 - \sqrt{3}) \ and \ D = 8\sqrt{2}/(2 - \sqrt{3})$.

8.3 Proof of Theorem 6.3

In order to apply the previous theorem, we first require the following (see [1, Prop. 4.11] or [34])³:

Lemma 8.3. The $(k,l)^{\text{th}}$ local coherence of the matrix UH^* , where U is the Walsh-Hadamard transform and H is the discrete Haar transform, satisfies $\mu_{k,k} = 1/(N_k - N_{k-1})$ and $\mu_{k,l} = 0$ for $k \neq l$.

Proof of Theorem 6.3. Fix $1 \le k \le r$, and let **s** be a sparsity pattern where $s_k = 0$ and $s_j \ge 1$ for at least one $j \ne k$. We compute a lower bound on the number of such pattern at the end of the proof.

We begin with (i) and choose \mathcal{M}_1 as a countably infinite subset of $H^*(\Sigma_{\mathbf{s},\mathbf{M}})$, where \mathbf{s} is as above. Let Ω be an (\mathbf{m}, \mathbf{N}) -multilevel random sampling scheme satisfying (6.3). Combining Theorem 8.2, Lemma 8.3 and (6.3), we see that

$$d^{H}(\Phi(Ax), x) \lesssim (1 + r^{1/4}) \frac{\sigma_{\mathbf{s}, \mathbf{M}}(Hx)_{\ell_{\boldsymbol{\omega}}^{1}}}{\sqrt{rs}} \equiv 0, \qquad \forall x \in H^{*}(\Sigma_{\mathbf{s}, \mathbf{M}}), \tag{8.7}$$

with probability greater than $1-\nu$, and hence, in particular, $d^H(\Phi(Ax),x)=0, \forall x\in \mathcal{M}_1$, with the same probability bound. It follows that Φ is an optimal map. Now let $x\in H^*(\Sigma_{\mathbf{s},\mathbf{M}})$, y=Ax and $e\in\mathbb{C}^m$ with $0<\|e\|_2\leqslant\varepsilon$. Recall the definition of Φ form (6.2), and let $\eta:=\inf\{\|y+e-z\|_{\ell^2}:z\in\mathcal{M}_2\}$. Let $x'\in H^*(\Sigma_{\mathbf{s},\mathbf{M}})$ be such that $\|y+e-Ax'\|_{\ell^2}=\eta$. Then, applying Theorem 8.2 with x' and y'=Ax'+e', where e'=y+e-Ax' satisfies $\|e'\|_{\ell^2}=\eta$, gives

$$d^{H}(\Phi(y+e), x') = d^{H}(\Phi(y'), x') \leqslant (1 + r^{1/4})D\eta.$$

Hence

$$d^{H}(\Phi(y+e), \Phi(y)) = d^{H}(\Phi(y+e), x) \leqslant ||x - x'||_{\ell^{2}} + (1 + r^{1/4})D\eta.$$
(8.8)

As noted, Theorem 8.2 is a simplified version of [1, Cor. 3.7]. Therein, it is shown that the matrix AH^* satisfies the so-called Restricted Isometry Property in levels (RIPL) of order (\mathbf{t}, \mathbf{M}) with $\delta_{\mathbf{t}, \mathbf{M}} \leq 1/2$, where

$$t_l = \min\{M_l - M_{l-1}, 2\lceil 4rs_l \rceil\} \geqslant \min\{M_l - M_{l-1}, 2s_l \}.$$

By definition, this means that $||z||_{\ell^2}^2/2 \le ||AH^*z||_{\ell^2}^2 \le 3||z||_{\ell^2}^2/2$ for all $z \in \Sigma_{\mathbf{t},\mathbf{M}}$. Notice that $H(x-x') \in \Sigma_{\mathbf{t},\mathbf{M}}$ since the difference of two (\mathbf{s},\mathbf{M}) -sparse vectors is at most $(2\mathbf{s},\mathbf{M})$ -sparse. It follows that

$$\begin{aligned} \|x - x'\|_{\ell^2} &= \|H(x - x')\|_{\ell^2} \leqslant \sqrt{2} \|A(x - x')\|_{\ell^2} \\ &= \sqrt{2} \|y - Ax'\|_{\ell^2} \leqslant \sqrt{2} (\|e\|_{\ell^2} + \eta) \leqslant 2\sqrt{2} \|e\|_{\ell^2}. \end{aligned}$$

Here, in the last step, we use the fact that $\eta \leq ||e||_{\ell^2}$, which follows from the definition of η . Returning to (8.8) and using this once more, we get

$$\frac{d^H(\Phi(y+e),\Phi(y))}{\|e\|_{\ell^2}} \leqslant 2\sqrt{2} + (1+r^{1/4})D.$$

³Note there is a minor difference between this and Prop. 4.11 of [1], in that in the setup above the first level includes only the scaling function, not both the scaling function and mother wavelet as was done in [1]. This difference is of no consequence to what follows.

Since $x \in H^*(\Sigma_{\mathbf{s},\mathbf{M}})$ and $0 < ||e||_{\ell^2} \leqslant \varepsilon$ were arbitrary, we obtain

$$L^{\varepsilon}(\Phi, \mathcal{M}_2) \leqslant L^{\varepsilon}(\Phi, AH^*(\Sigma_{\mathbf{s}, \mathbf{M}})) \leqslant 2\sqrt{2} + (1 + r^{1/4})D,$$

as required.

We now consider (ii). Let $c \in \mathbb{C}^N$ satisfy $||c||_{\ell^2} = 1$ and $c_i = 0$ for $i \notin \{N_{k-1} + 1, \dots, N_k\}$. Define $z_1 = H^*c$ and let $0 \neq z_2 \in \mathcal{M}_1$. Let $z = z_1 + \frac{\kappa}{\|z_2\|_{\ell^2}} z_2$, where $\kappa \in \mathbb{R}$ is a parameter to be set later, and define

$$\widetilde{\mathcal{M}}_1 = \mathcal{M}_1 \cup \{0, z\}.$$

Since $m_k = 0$ and the matrix UH^* is block diagonal (see Lemma 8.3), it follows that $z_1 \in \text{Ker}(A)$. Hence, by (8.7), the fact that $||z_1||_{\ell^2} = 1$ (as the Haar transform is unitary), and the assumption (6.4) on Ψ , we have that

$$\begin{split} \frac{1}{p} \sup_{x \in \widetilde{\mathcal{M}}_1} d^H(\Phi(Ax), x) &\geqslant \frac{1}{p} d^H(\Phi(Az), z) = \frac{1}{p} d^H(\Phi(A\frac{\kappa}{\|z_2\|_{\ell^2}} z_2), z) \\ &= \frac{1}{p} \|\frac{\kappa}{\|z_2\|_{\ell^2}} z_2 - z\|_{\ell^2} = \frac{1}{p} \|z_1\|_{\ell^2} \geqslant \delta \geqslant \sup_{x \in \widetilde{\mathcal{M}}_1} d^H(\Psi(Ax), x), \end{split}$$

which establishes (6.5). Now consider (6.6). As shown in [1], the conditions on **m** assert that A satisfies the so-called Restricted Isometry Property in Levels (see Definition 3.4 therein). Hence $||Az||_{\ell^2} = \frac{\kappa}{||z_2||_{\ell^2}}||Az_2||_{\ell^2} \leq 2\kappa$. Therefore, if $2\kappa \leq \varepsilon$, we have

$$L^{\varepsilon}(\Psi,\widetilde{\mathcal{M}}_2)\geqslant L^{\varepsilon}(\Psi,0)\geqslant \frac{d^H(\Psi(Az),\Psi(0))}{\|Az\|_{\ell^2}}\geqslant \frac{\|z\|_{\ell^2}-2\delta}{2\kappa}\geqslant \frac{1-2/p}{2\kappa}.$$

Setting $\kappa = (1 - 2/p)/(2p) \leqslant \varepsilon/2$ now gives $L^{\varepsilon}(\Psi, \widetilde{\mathcal{M}}_2) \geqslant p \geqslant 1/\varepsilon$. On the other hand, for the mapping Φ notice that

$$L^{\varepsilon}(\Phi,\widetilde{\mathcal{M}}_{2}) = L^{\varepsilon}(\Phi,\mathcal{M}_{2}) \vee L^{\varepsilon}(\Phi,0) \vee L^{\varepsilon}(\Phi,Az)$$

Notice that $L^{\varepsilon}(\Phi, \mathcal{M}_2) \vee L^{\varepsilon}(\Phi, 0) \leqslant L^{\varepsilon}(\Phi, AH^*(\Sigma_{\mathbf{s}, \mathbf{M}}))$. Also, since $Az = A(\frac{\kappa}{\|z_2\|_{\ell^2}}z_2)$, we have

$$\begin{split} L^{\varepsilon}(\Phi,Az) &= \sup_{0 < \|e\|_{\ell^{2}} \leqslant \varepsilon} \frac{d^{H}(\Phi(Az+e),\Phi(Az))}{\|e\|_{\ell^{2}}} \\ &= L^{\varepsilon}\left(\Phi,A\left(\frac{\kappa}{\|z_{2}\|_{\ell^{2}}}z_{2}\right)\right) \leqslant L^{\varepsilon}(\Phi,AH^{*}(\Sigma_{\mathbf{s},\mathbf{M}})). \end{split}$$

Hence, by this and the proof of part (i), we have

$$L^{\varepsilon}(\Phi, \widetilde{\mathcal{M}}_2) \leqslant L^{\varepsilon}(\Phi, AH^*(\Sigma_{\mathbf{s}, \mathbf{M}})) \leqslant 2\sqrt{2} + (1 + r^{1/4})D$$

as required. The statement about uncountably many domains $\mathcal{M}_1, \widetilde{\mathcal{M}}_1$ follows by multiplying the above choices by complex numbers of modulus one.

To see that there will be at least $2^{(r-1)(r-2)/2}$ sparsity patterns where $s_k=0$, for some $1 \le k \le r$, and at least one $s_j>0$ for some $j \ne k$, consider the case where $s_1=0$ and $s_j \ge 1$ for j>1. As $\mathbf{M}=(2^0,2^1,2^2,\ldots,2^{r-1})$ the size of the r-1 last levels are $2^k,k=0,\ldots,r-2$. Thus there are

$$\prod_{k=0}^{r-2} 2^k = 2^{\sum_{k=1}^{r-2} k} = 2^{(r-1)(r-2)/2}$$

sparsity patterns, where $s_1 = 0$ and $s_j \ge 1$ for j > 1.

8.4 Proof of Theorem 6.5

We first require the following result. This is based on [3, Lem. 1], although we slightly improve the scaling for the case l > k (the proof is identical):

Lemma 8.4. The (k,l)th local coherence of the matrix UH^* , where U is the discrete Fourier transform and H is the discrete Haar transform, satisfies

$$\mu_{k,l} \lesssim \frac{1}{N_k - N_{k-1}} \left\{ \begin{array}{ll} 2^{-(k-l)} & l \leqslant k \\ 2^{-3(l-k)} & l > k \end{array} \right., \qquad k, l = 1, \dots, r.$$

Proof of Theorem 6.5. Consider part (i). Pick \mathcal{M}_1 as a subset of $H^*(\Sigma_{\mathbf{s},\mathbf{M}})$ of size M and let Ω be an (\mathbf{m},\mathbf{N}) -multilevel random sampling scheme satisfying (6.8). The remainder of the proof is now identical to that of part (i) of Theorem 6.3. The only difference is the us of Lemma 8.4 instead of Lemma 8.3.

Consider part (ii). Since \mathcal{NN} is M-interpolatory and $|\mathcal{M}_1| = M$ the first result follows immediately. Consider the second result. Let $x_{K+i} = x_i + z_i$. Since $m \leqslant N-1$ we may choose $z_2 \in \mathcal{N}(A)$ with $||z_2|| = \gamma/4$. We now claim there exists z_1 such that $||z_1|| = \gamma$ and $0 < ||Az_1|| \leqslant \gamma^2/2$. Indeed, let $u \in \mathcal{N}(A)$ and $v \in \mathcal{N}(A)^{\perp}$ be unit vectors and define $z_1 = \sqrt{\gamma^2 - \beta^2}u + \beta v$, where $\beta = \gamma^2/(2||A||)$. Notice that $||A|| \geqslant 1$, since $A = P_{\Omega}DF$, where F is unitary and D is diagonal with diagonal entries at least one, and therefore $\beta \leqslant \gamma^2/2 < \gamma$. Hence z_1 is well defined. The claim now follows immediately.

Since \mathcal{NN} is (K+1)-interpolatory, there exists a network $\widetilde{\Psi} \in \mathcal{NN}$ with $\widetilde{\Psi}(Ax_i) = x_i$ for $i = 1, \ldots, K+1$. Notice that $Ax_{K+2} = Ax_2$. Hence

$$||x_i - \widehat{\Psi}(Ax_i)||_{\ell^2} \le ||x_{K+2} - \widetilde{\Psi}(Ax_{K+2})||_{\ell^2} = ||x_{K+2} - x_2||_{\ell^2}, \quad i = 1, \dots, K+2.$$

Hence, since $0 < ||Az_1|| \le \gamma^2/2 \le \varepsilon$, we have

$$L^{\varepsilon}(\widehat{\Psi}, \mathcal{M}_{2}) \geqslant \frac{\|\widehat{\Psi}(Ax_{1} + Az_{1}) - \widehat{\Psi}(Ax_{1})\|_{\ell^{2}}}{\|Az_{1}\|_{\ell^{2}}}$$

$$\geqslant \frac{\|\widehat{\Psi}(Ax_{K+1}) - \widehat{\Psi}(Ax_{1})\|_{\ell^{2}}}{\gamma^{2}/2}$$

$$\geqslant \frac{\|x_{K+1} - x_{1}\|_{\ell^{2}} - \|\widehat{\Psi}(Ax_{K+1}) - x_{K+1}\|_{\ell^{2}} - \|\widehat{\Psi}(Ax_{1}) - x_{1}\|_{\ell^{2}}}{\gamma^{2}/2}$$

$$\geqslant \frac{\|x_{K+1} - x_{1}\|_{\ell^{2}} - 2\|x_{K+2} - x_{2}\|_{\ell^{2}}}{\gamma^{2}/2}$$

$$= \frac{\|z_{1}\|_{\ell^{2}} - 2\|z_{2}\|_{\ell^{2}}}{\gamma^{2}/2} = \frac{\gamma/2}{\gamma^{2}/2} = \frac{1}{\gamma},$$

as required.

8.5 Proof of Theorem 7.2

Proof Theorem 7.2. Let $\widetilde{\Omega} \subset \{1, \dots, N\}$ such that $2 \leqslant |\widetilde{\Omega}| \leqslant N-2$. Note that there are $2^N - 2N - 2$ different choices of $\widetilde{\Omega}$.

Moreover, let $\Omega = \widetilde{\Omega} \cup \{j\}$ where $j \notin \widetilde{\Omega}$, so $\Omega \neq \{1, ..., N\}$ by the fact that $|\widetilde{\Omega}| \leqslant N - 2$. Choose $\mathcal{M}_1 \subset \mathcal{N}(\widetilde{A})^{\perp}$ of size K. If K is infinite, choose \mathcal{M}_1 to be countable. By multiplying \mathcal{M}_1 by any real number we clearly get uncountably many different choices of \mathcal{M}_1 . For ease of notation let $\mathcal{M}_2 = A(\mathcal{M}_1)$ and $\widetilde{\mathcal{M}}_2 = \widetilde{A}(\mathcal{M}_1)$.

Choose $i \in \widetilde{\Omega}$. Let $\widetilde{\mathcal{T}} \subset \widetilde{\mathcal{M}}_2 \times \mathcal{M}_1$ be any finite non-zero collection such that there is exactly one pair $(\widehat{y} = \widetilde{A}\widehat{x}, \widehat{x}) \in \widetilde{\mathcal{T}}$ such that

$$\widetilde{\mathcal{T}} \setminus \{(\widehat{y}, \widehat{x})\} \subset \{(\widetilde{y}, \widetilde{x}) \mid \widetilde{y} = \widetilde{A}\widetilde{x}, P_i \widetilde{y} = 0\}, \quad P_i \widehat{y} \neq 0, \quad P_{\widetilde{\Omega} \setminus \{i\}} \widehat{y} = 0, \tag{8.9}$$

where P_i denotes the projection onto the *i*-th coordinate. Note that such a choice is possible since $2 \leq |\widetilde{\Omega}|$. Choose any non-empty $\mathcal{T} \subset \mathcal{M}_2 \times \mathcal{M}_1$. As both \mathcal{T} and $\widetilde{\mathcal{T}}$ can be multiplied by any real number that would not change any of the properties outlined above, we clearly have uncountably many different choices of \mathcal{T} and $\widetilde{\mathcal{T}}$. Note that

$$\widetilde{A}^{\dagger} y = x \text{ if } y = \widetilde{A}x, \quad x \in \mathcal{M}_1,$$
 (8.10)

where \widetilde{A}^{\dagger} denotes the pseudoinverse. This fact will be crucial later in the argument.

Next we write all networks as complex valued for clarity, yet it should implicitly be understood that they can be written as real-valued, by doubling all dimensions (see Remark 2.1). Consider the L-layer ReLU neural network $\widetilde{\Psi}: \mathbb{C}^{|\widetilde{\Omega}|} \to \mathbb{C}^N$ defined by

$$\widetilde{\Psi}(x) = \widetilde{A}^{\dagger} W_2 \rho(\dots \rho(W_1 W_2 \rho(W_1 W_2 \rho(W_1 x)))),$$

$$W_1 = [1, -1]^T \otimes I_{|\widetilde{\Omega}|}, \quad W_2 = [1, -1] \otimes I_{|\widetilde{\Omega}|},$$
 (8.11)

where we use the notation I_d for the d-dimensional identity matrix and \otimes denotes the Kronecker product.

Observe that for any pair $(\widetilde{y} = \widetilde{A}\widetilde{x}, \widetilde{x}) \in \widetilde{\mathcal{M}}_2 \times \mathcal{M}_1$ we have

$$\widetilde{\Psi}(\widetilde{y}) = \widetilde{A}^{\dagger} W_{2} \rho(\dots \rho(W_{1} W_{2} \rho(W_{1} W_{2} \rho(W_{1} \widetilde{y}))))
= \widetilde{A}^{\dagger} W_{2} \rho(\dots \rho(W_{1} W_{2} \rho(W_{1}([1, -1] \otimes I_{|\widetilde{\Omega}|}) \rho(([1, -1]^{T} \otimes I_{|\widetilde{\Omega}|}) \widetilde{y}))))
= \widetilde{A}^{\dagger} W_{2} \rho(\dots \rho(W_{1} W_{2} \rho(W_{1}([1, -1] \otimes I_{|\widetilde{\Omega}|}) [\rho(\widetilde{y}), \rho(-\widetilde{y})]^{T})))
= \widetilde{A}^{\dagger} W_{2} \rho(\dots \rho(W_{1} W_{2} \rho(W_{1}(\rho(\widetilde{y}) - \rho(-\widetilde{y}))))) = \widetilde{A}^{\dagger} \widetilde{y} = \widetilde{x},$$
(8.12)

which follows from (8.11), (8.10) and the easy observation that, since ρ is the ReLU function, $\rho(\widetilde{y}) - \rho(-\widetilde{y}) = \widetilde{y}$. Also, recalling that $\mathbf{n} = (m, n_1, \dots, n_{L-1}, N)$, $m = \max\{|\Omega|, |\widetilde{\Omega}|\}$ and the assumption that $n_j \geq 2m$, it is clear that by replacing W_j with $W_j \oplus 0$ we can without loss of generality assume that $\widetilde{\Psi} \in \mathcal{NN}_{\mathbf{n}}$. Hence, by setting λ_{opt} to zero, and using (8.12) we have that

$$\sum_{(\widetilde{y},\widetilde{x})\in\widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \widetilde{\Psi}(\widetilde{y})\|_{\ell^{2}}^{2} + \lambda_{\text{opt}} J(\widehat{\Psi}) = 0.$$
(8.13)

We will continue with this choice of λ_{opt} throughout the argument. We claim that λ_{opt} is optimal for $\{(\widetilde{A}, \mathcal{M}_1), \mathcal{NN}, \widetilde{\mathcal{T}}\}$. Indeed, (8.13) implies that $\widetilde{\Psi}$ is a minimiser of (7.4) for $\lambda = \lambda_{\text{opt}}$. Thus, we only need to show that $\widetilde{\Psi}$ is an optimal mapping. To see this we observe that by (8.12) it follows that

$$c_{\text{opt}}(\widetilde{A}, \mathcal{M}_{1}) = \inf_{\varphi : \widetilde{\mathcal{M}}_{2} \rightrightarrows \mathcal{M}_{1}} \sup_{x \in \mathcal{M}_{1}} d_{\mathcal{M}_{1}}^{H}(\varphi(\widetilde{A}x), x)$$

$$\leq \sup_{x \in \mathcal{M}_{1}} d_{\mathcal{M}_{1}}^{H}(\widetilde{\Psi}(\widetilde{A}x), x) = 0,$$
(8.14)

and hence our claim that λ_{opt} is optimal for $\{(\widetilde{A}, \mathcal{M}_1), \mathcal{NN}, \widetilde{\mathcal{T}}\}$ is true. That λ_{opt} also is optimal for $\{(A, \mathcal{M}_1), \mathcal{NN}, \mathcal{T}\}$, will be shown when we consider part (ii).

We will now establish the collection S as described in the statement of the theorem. Note that $\mathcal{N}(\widetilde{A}) \cap \mathcal{N}(A)^{\perp}$, is non-zero by our choice of Ω . We choose a non-zero $x \in \mathcal{N}(\widetilde{A}) \cap \mathcal{N}(A)^{\perp}$ and let $z = \widehat{x} + x$ where \widehat{x} is from (8.9), and let S denote any non-zero uncountable collection of multiples of z.

To show (i), we note that it is enough to show that λ_{opt} is no longer optimal if we add or replace a specific element of $\widetilde{\mathcal{T}}$ with an element from \mathcal{S} . First, suppose we replace $(\widehat{y}, \widehat{x}) \in \widetilde{\mathcal{T}}$ by $(\widehat{y}, \widehat{x} + x)$, and define the neural network $\Phi \in \mathcal{NN}$ by

$$\Phi(x) = T\rho(\dots \rho(W_1 W_2 \rho(W_1 W_2 \rho(W_1 x)))), \tag{8.15}$$

where W_1 and W_2 are defined in (8.11) and

$$T = C([1,-1] \otimes I_{|\widetilde{\Omega}|}), \qquad C = \widetilde{A}^{\dagger} P_{\widetilde{\Omega} \backslash \{i\}} + (\frac{1}{\widehat{y}^i} (\widehat{x} + x) \otimes e_i^T).$$

We observe that for $(\widetilde{y}, \widetilde{x}) \in \widetilde{\mathcal{T}}$ with $P_i \widetilde{y} = 0$ we have $\Phi(\widetilde{y}) = \widetilde{\Psi}(\widetilde{y}) = \widetilde{x}$, and for $(\widehat{y}, \widehat{x} + x) \in \widetilde{\mathcal{T}}$ we have by (8.9) that

$$\Phi(\widehat{y}) = \left(\widetilde{A}^{\dagger} P_{\widetilde{\Omega} \setminus \{i\}} + \left(\frac{1}{\widehat{y}^{i}}(\widehat{x} + x) \otimes e_{i}^{T}\right)\right) \widehat{y} = \widehat{x} + x.$$

Hence, the objective function in (7.4) is zero at Φ whenever $\lambda = \lambda_{\rm opt}$. Thus, any

$$\widehat{\Phi} \in \underset{\widehat{\Psi} \in \mathcal{NN}}{\operatorname{argmin}} \frac{1}{|\widetilde{\mathcal{T}}|} \sum_{(\widetilde{y}, \widetilde{x}) \in \widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \widehat{\Psi}(\widetilde{y})\|_{\ell^{2}}^{2} + \lambda_{\operatorname{opt}} J(\widehat{\Psi})$$
(8.16)

will satisfy $\widehat{\Phi}(\widehat{y}) = \widehat{x} + x$, which means that $\sup_{x \in \mathcal{M}_1} d_{\mathcal{M}_1}^H(\widehat{\Phi}(\widetilde{A}x), x) \neq 0$ which by (8.14) means that $\widehat{\Phi}$ is not an optimal map. Hence λ_{opt} is no longer optimal for $\{(\widetilde{A}, \mathcal{M}_1, \mathcal{NN}, \widetilde{\mathcal{T}})\}$.

Now let us consider the case where z is added to $\widetilde{\mathcal{T}}$. Consider the neural net $\widetilde{\Phi}$ defined by (8.15) with

$$T = D([1, -1] \otimes I_{|\widetilde{\Omega}|}), \qquad D = \widetilde{A}^{\dagger} P_{\widetilde{\Omega} \setminus \{i\}} + (\frac{1}{\widehat{y}^{i}} (\widehat{x} + \frac{1}{2} x) \otimes e_{i}^{T}).$$

For $(\widetilde{y}, \widetilde{x}) \in \widetilde{\mathcal{T}}$ with $P_i \widetilde{y} = 0$ we have $\widetilde{\Phi}(\widetilde{y}) = \widetilde{x}$, and for $(\widehat{y}, \widehat{x} + x), (\widehat{y}, \widehat{x}) \in \widetilde{\mathcal{T}}$ we have by (8.9) that

$$\widetilde{\Phi}(\widehat{y}) = \left(\widetilde{A}^{\dagger} P_{\widetilde{\Omega} \setminus \{i\}} + (\frac{1}{\widehat{y}^i} (\widehat{x} + \frac{1}{2} x) \otimes e_i^T)\right) \widehat{y} = \widehat{x} + \frac{1}{2} x.$$

Hence, $\frac{1}{|\widetilde{\mathcal{T}}|} \sum_{(\widetilde{y},\widetilde{x}) \in \widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \widetilde{\Phi}(\widetilde{y})\|_{\ell^2}^2 + \lambda_{\text{opt}} J(\widetilde{\Phi}) = \frac{1}{|\widetilde{\mathcal{T}}|} \|\frac{1}{2}x\|_{\ell^2}^2$, and therefore any minimiser $\widehat{\Psi}$ of (8.16) will satisfy

$$\frac{1}{|\widetilde{\mathcal{T}}|} \sum_{(\widetilde{y},\widetilde{x}) \in \widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \widehat{\Psi}(\widetilde{y})\|_{\ell^{2}}^{2} + \lambda_{\mathrm{opt}} J(\widehat{\Psi}) \leqslant \frac{1}{4|\widetilde{\mathcal{T}}|} \|x\|_{\ell^{2}}^{2}.$$
(8.17)

However, by (8.14), any optimal map φ for $(\widetilde{A}, \mathcal{M}_1)$ will satisfy

$$\sup_{x \in \mathcal{M}_1} d_{\mathcal{M}_1}^H(\varphi(\widetilde{A}x), x) = 0.$$

Thus, $\frac{1}{|\widetilde{\mathcal{T}}|} \sum_{(\widetilde{y},\widetilde{x}) \in \widetilde{\mathcal{T}}} \frac{1}{2} \|\widetilde{x} - \varphi(\widetilde{y})\|_{\ell^2}^2 = \frac{1}{2|\widetilde{\mathcal{T}}|} \|x\|_{\ell^2}^2$, and therefore, by (8.17), no minimiser $\widehat{\Psi}$ of (8.16) can be an optimal map for $(\widetilde{A}, \mathcal{M}_1)$. Hence λ_{opt} is not optimal for $\{(\widetilde{A}, \mathcal{M}_1, \mathcal{NN}, \widetilde{\mathcal{T}})\}$.

Consider part (ii). Let $\mathcal{D} \subset \mathcal{M}_1 \cup \mathcal{S}$ be any finite non-empty set, and let $\mathcal{T} = \{(Ax, x) : x \in \mathcal{D}\}$. We shall prove that λ_{opt} is optimal for $\{(A, \mathcal{N}(A)^{\perp}), \mathcal{N}\mathcal{N}, \mathcal{T}\}$ for any such \mathcal{T} . Note that this is a stronger statement than in the theorem, as $\mathcal{M}_1 \cup \mathcal{S} \subset \mathcal{N}(A)^{\perp}$. From (8.14) it is clear that $y \mapsto A^{\dagger}y$ is an optimal map for $(A, \mathcal{N}(A)^{\perp})$. Using the network $\widetilde{\Psi}$ from (8.12), where \widetilde{A}^{\dagger} is replaced by A^{\dagger} in the last layer, it is clear that $\widetilde{\Psi}$ is an optimal map, and a minimiser of (8.16) for λ_{opt} when we sum over \mathcal{T} . Thus λ_{opt} is optimal for $\{(A, \mathcal{N}(A)^{\perp}), \mathcal{N}\mathcal{N}, \mathcal{T}\}$.

9 Methods

In this section, we describe the methods used to generate the various numerical results.

9.1 Sparse regularization decoders

The sparse regularization decoder used in Figures 5-8, is

$$\underset{z \in \mathbb{C}^N}{\text{minimise}} \|z\|_{\ell^1} \text{ subject to } \|AH^\top z - y\|_{\ell^2}^2 \leqslant \eta, \tag{9.1}$$

where $\eta \geqslant 0$ is a noise parameter and $H \in \mathbb{R}^{N \times N}$ is a discrete wavelet transform with a Daubechies orthonormal wavelet. The sampling operator A is a subsampled discrete Fourier transform in all of these experiments. To search for a minimiser of the above optimization problem, we used the SPGL1 [47] software package. The chosen parameters in each figure can be found in Table 1.

For Fig. 3 and Fig. 4, we used the sparse regularization decoder introduced in [30]. This is a more advanced decoder, supporting different types of regularizes and also non-Fourier sampling operators, such as Radon sampling.

The decoder proposed in [30] tries to iteratively solve the optimization problem

$$\underset{z \in \mathbb{C}}{\text{minimise}} \sum_{j=1}^{J} \lambda_j \|W_j H_j z\|_{\ell^1} + \text{TGV}_{\alpha}^2(z) \text{ subject to } Az = y$$

$$(9.2)$$

using n iterations. Here the W_j 's are diagonal weighting matrices, $\lambda_j \in \mathbb{R}_+$ are weighting parameters, and H_j is the j'th subband in a shearlet transform. The weights W_j and λ_j are updated iteratively between each iteration. The $\mathrm{TGV}_{\alpha}^2(z)$ term is a second order Total Generalised Variation operator depending on two parameters $\alpha = (\alpha_1, \alpha_2)$, where the first order term (TV) is weighted by α_1 and a second order (generalised) term weighted by α_2 .

In all experiments we used shearlets with 4 scales and directional parameters [0,0,1,1]. For the experiment in Fig. 3, we set $\alpha_2 = 0$, i.e. the TGV-term simplifies to TV-minimization. The complete set of parameters can be found in Table 1.

9.2 Training of the network in Fig. 5

In this section we will describe the network architecture and training procedure for the network trained to reconstruct lines in Fig. 5 from Fourier measurements. We do not go into details of architectures and training procedures of any of the other networks, as this information can be found in the corresponding papers. See also supplementary information of [6] for a summary.

For the experiment in Fig. 5, we used a ground truth dataset consisting of 75+6 images, that is vectors in \mathbb{R}^N with N=512. The 75 first images was created by evaluating piecewise continuous polynomials p(t) at an equidistant grid on [0,1].

Thus our ground truth data $x^{(k)} \in \mathbb{R}^N$ was created as $x_i^{(k)} = p(t_i)$, i = 1, ..., N, for grid points t_i . The piecewise continuous polynomials p(t) where all created by drawing coefficients $c_1 \sim \mathcal{U}([0,10))$, $b_1, s_i, a_i, \sim \mathcal{U}([0,1))$, i = 1, 2, 3 and $d_1, d_2, d_3 \sim \mathcal{U}([-1/2, 1/2))$, from uniform distributions and forming p as

$$p(t) := c_1(t - a_1)(z - a_2)(t - a_3) + b_1 + d_1h(t - s_1) + d_2h(t - s_2) + d_3h(t - s_3)$$

where $h(t) = \mathbb{1}_{t \ge 0}$ is the Heaviside step function.

| Parameters for (9.1) | | | Paramet | Parameters for (9.2) | | | |
|----------------------|---------|--------|---------|------------------------|------------|------------|--|
| | Wavelet | η | | | | | |
| Fig. 5. | DB4 | 0.01 | | n | α_1 | α_2 | |
| Fig. 6. | DB2 | 0.1 | Fig. 3. | 50 | 1 | 0 | |
| Fig. 7. | DB2 | 0.1 | Fig. 4. | 500 | 1 | 1 | |
| Fig. 8. | DB4 | 0.001 | | | | | |

Table 1: Parameters for the sparse regularization decoders.

We formed the last six training images by picking the vectors $x^{(k)} \in \mathbb{R}^N$, $k=1,\ldots,6$ and add the vectors $x^{(k+75)} = x^{(k)} + z$, $k=1,\ldots,6$ to the ground truth dataset. Here $z \in \mathbb{R}^N$ is the high frequent "Tumor" displayed in Fig. 5. Thus the complete training set was formed as

$$\mathcal{T} = \{ (Ax^{(k)}, x^{(k)})_{k=1}^{81} \subset \mathbb{R}^N \times \mathbb{C}^m \}$$

where the sampling matrix $A = P_{\Omega}F \in \mathbb{C}^{m \times N}$, m = 77 is a subsampled discrete Fourier transform

For the neural network $f: \mathbb{C}^m \to \mathbb{R}^N$ we choose a U-net [39] like architecture, inspired by the FBPConvNet proposed in [28]. That is, we let $f(y) = \varphi(A^*y)$ where $\varphi: \mathbb{R}^{N \times 2} \to \mathbb{R}^N$ is a neural network with the U-net like structure. Here the input $A^*y \in \mathbb{C}^N$ to φ is split into its real and imaginary component and processed as an element of $\mathbb{R}^{N \times 2}$. The complete architecture of φ can be seen in Fig. 9. The network f was trained for 100,000 epochs using the Adam optimizer, with a batch size of 75. The dataset was shuffled between each epoch.

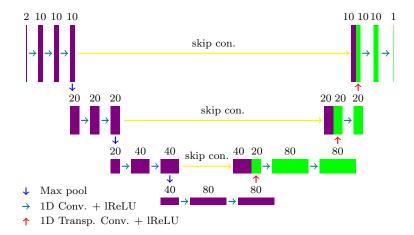


Figure 9: The network architecture of of φ . The numbers above all the squares indicate the number of channels of the tensor. The arrows indicate different types of tensor operations. The teal arrows symmetrically extend the tensor along the resolution axis, so that after preforming a 1D convolution (on the symmetrically extended signal) with kernel size 3 and "vaild" boundary handling, the resolution remained the same. The convolution is followed by a leaky ReLU (lReLU) activation function with parameter $\alpha=0.2$. The blue arrow indicate a max pooling of size 2 in the resolution axis. The red arrow is a transpose 1D convolution with kernel size 3 and "same" boundary handling, followed by a lReLU activation function ($\alpha=0.2$). Finally the yellow arrows are skip connections, concatenating an earlier tensor with the output from a different layer in the network.

9.3 Creating Gaussian noise in $\mathcal{N}(A)^{\perp}$

In Fig. 4 and 6 we construct Gaussian vectors $v \in \mathbb{C}^N$ of a fixed magnitude, all lying in $\mathcal{N}(A)^{\perp}$. This is done as follows. We draw the real and imaginary components of a vector $e \in \mathbb{C}^m$ form a Normal distribution $\mathcal{N}(0,10)$ and compute v as $\alpha A^*e = v$, where $\alpha \in \mathbb{R}_+$ is a scalar. The α is chosen so that v gets the desired norm. A Gaussian random variable is still Gaussian after a linear map, so the vector v is Gaussian. Furthermore, since since $AA^* = I$ for the special case where $A = P_{\Omega}F$ is a subsampled discrete Fourier transform, it follows that $v \in \mathcal{N}(A)^{\perp}$.

References

[1] B. Adcock, V. Antun, and A. C. Hansen. Uniform recovery in infinite-dimensional compressed sensing and applications to structured binary sampling. arXiv:1905.00126, 2019.

- [2] B. Adcock, A. C. Hansen, C. Poon, and B. Roman. Breaking the coherence barrier: A new theory for compressed sensing. In *Forum of Mathematics, Sigma*, volume 5. Cambridge University Press, 2017.
- [3] B. Adcock, A. C. Hansen, and B. Roman. A note on compressed sensing of structured sparse wavelet coefficients from subsampled Fourier measurements. arXiv:1403.6541, 2014.
- [4] B. Adcock, A. C. Hansen, and B. Roman. The quest for optimal sampling: Computationally efficient, structure-exploiting measurements for compressed sensing. In *Compressed Sensing and Its Applications*. Birkhäuser, 2015.
- [5] J. Adler and O. Öktem. Solving ill-posed inverse problems using iterative deep neural networks. *Inverse Problems*, 33:124007, 2017.
- [6] V. Antun, F. Renna, C. Poon, B. Adcock, and A. C. Hansen. On instabilities of deep learning in image reconstruction does AI come at a cost? arXiv preprint arXiv:1902.05300, 2019.
- [7] S. Arridge, P. Maass, O. Öktem, and C.-B. Schönlieb. Solving inverse problems using data-driven models. *Acta Numer.*, 28:1–174, 2019.
- [8] N. Baker, F. Alexander, T. Bremer, A. Hagberg, Y. Y. Kevrekidis, H. Najm, M. Parashar, A. Patra, J. Sethian, S. Wild, and K. Willcox. Workshop report on basic research needs for scientific machine learning: Core technologies for artificial intelligence. U.S. Department of Energy Advanced Scientific Computing Research, 2019.
- [9] A. Bastounis and A. C. Hansen. On the absence of uniform recovery in many real-world applications of compressed sensing and the restricted isometry property and nullspace property in levels. SIAM J. Imaging Sci., 10(1):335–371, 2017.
- [10] T. A. Bubba, G. Kutyniok, M. Lassas, M. März, W. Samek, S. Siltanen, and V. Srinivasan. Learning the invisibile: A hybrid deep learning-shearlet framework for limited angle computed tomography. arXiv:1811.04602, 2018.
- [11] E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, 52(2):489–509, 2006.
- [12] E. J. Candès and T. Tao. Near optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. Inform. Theory*, 52(12):5406–5425, 2006.
- [13] A. Chakraborty, M. Alam, V. Dey, A. Chattopadhyay, and D. Mukhopadhyay. Adversarial attacks and defences: A survey. arXiv:1810.00069, 2018.
- [14] A. Cohen, W. Dahmen, and R. A. DeVore. Compressed sensing and best k-term approximation. J. Amer. Math. Soc., 22(1):211–231, 2009.
- [15] M. Dashti and A. M. Stuart. The bayesian approach to inverse problems. In *Handbook of Uncertainty Quantification*. Springer, 2017.
- [16] D. L. Donoho. Compressed sensing. IEEE Trans. Inform. Theory, 52(4):1289–1306, 2006.
- [17] W. E, J. Han, and A. Jentzen. Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. *Com*mun. Math. Stat., 5(4):349–380, 2017.
- [18] Q. Fan, T. Witzel, A. Nummenmaa, K. R. Van Dijk, J. D. Van Horn, M. K. Drews, L. H. Somerville, M. A. Sheridan, R. M. Santillana, J. Snyder, et al. MGH–USC human connectome project datasets with ultra-high b-value diffusion MRI. *Neuroimage*, 124:1108– 1114, 2016.
- [19] FDA. FDA permits marketing of artificial intelligence-based device to detect certain diabetes-related eye problems. FDA News Release, April 11 2018, 2018.

- [20] S. G. Finlayson, J. D. Bowers, J. Ito, J. L. Zittrain, A. L. Beam, and I. S. Kohane. Adversarial attacks on medical machine learning. *Science*, 363(6433):1287–1289, 2019.
- [21] S. Foucart and H. Rauhut. A mathematical introduction to compressive sensing. Birkhauser, 2013.
- [22] I. J. Goodfellow, Y. Bengio, and A. Courville. Deep learning. MIT Press, 2016.
- [23] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In Proceedings of the Int. Conf. on Learning Representations, 2015.
- [24] H. Gupta, K. H. Jin, H. Q. Nguyen, M. T. McCann, and M. Unser. CNN-based projected gradient descent for consistent CT image reconstruction. *IEEE transactions on medical* imaging, 37(6):1440–1453, 2018.
- [25] C. F. Higham and D. J. Higham. Deep learning: An introduction for applied mathematicians. SIAM Rev., 61(4):860–891, 2019.
- [26] Y. Huang, A. Preuhs, G. Lauritsch, M. Manhart, X. Huang, and A. Maier. Data consistent artifact reduction for limited angle tomography with deep learning prior. In *Int. Works. on Machine Learning for Medical Image Reconstruction*, pages 101–112. Springer, 2019.
- [27] Y. Huang, T. Würfl, K. Breininger, L. Liu, G. Lauritsch, and A. Maier. Some investigations on robustness of deep learning in limited angle tomography. In *Int. Conf. on Medical Image Computing and Computer-Assisted Intervention*, pages 145–153. Springer, 2018.
- [28] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser. Deep convolutional neural network for inverse problems in imaging. *IEEE Trans. Image Process.*, 26(9):4509–4522, 2017.
- [29] C. Li and B. Adcock. Compressed sensing with local structure: Uniform recovery guarantees for the sparsity in levels class. *Appl. Comput. Harmon. Anal.*, 2017.
- [30] J. Ma and M. März. A multilevel based reweighting algorithm with joint regularizers for sparse recovery. arXiv preprint arXiv:1604.06941, 2016.
- [31] M. T. McCann, K. H. Jin, and M. Unser. Convolutional neural networks for inverse problems in imaging: A review. *IEEE Signal Process. Mag.*, 34(6):85–95, 11 2017.
- [32] S. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard. Universal adversarial perturbations. In *IEEE Conf. on computer vision and pattern recognition*, pages 86–94, July 2017.
- [33] S. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *IEEE Conf. on Computer Vision and Pattern Recognition*, pages 2574–2582, June 2016.
- [34] A. Moshtaghpour, J. M. B. Dias, and L. Jacques. Close encounters of the binary kind: Signal reconstruction guarantees for compressive Hadamard sampling with Haar wavelet basis. arXiv preprint arXiv:1907.09795, 2019.
- [35] A. Pinkus. Approximation theory of the MLP model in neural networks. *Acta numer.*, 8:143–195, 1999.
- [36] C. Poon. On the role of total variation in compressed sensing. SIAM J. Imaging Sci., 8(1):682–720, 2015.
- [37] B. Roman, A. Bastounis, B. Adcock, and A. C. Hansen. On fundamentals of models and sampling in compressed sensing. *Preprint*, 2015.
- [38] B. Roman, A. C. Hansen, and B. Adcock. On asymptotic structure in compressed sensing. arXiv:1406.4178, 2014.

- [39] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *Int. Conf. on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [40] S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4), 2017.
- [41] J. Schlemper, J. Caballero, J. V. Hajnal, A. Price, and D. Rueckert. A deep cascade of convolutional neural networks for MR image reconstruction. In *Int. Conf. on Information Processing in Medical Imaging*, pages 647–658. Springer, 2017.
- [42] C. Schwab and J. Zech. Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ. *Analysis and Applications*, 17(01):19–55, 2019.
- [43] A. Sriram, J. Zbontar, T. Murrell, C. L. Zitnick, A. Defazio, and D. K. Sodickson. GrappaNet: Combining parallel imaging with deep learning for multi-coil MRI reconstruction. arXiv preprint arXiv:1910.12325, 2019.
- [44] R. Strack. Imaging: AI transforms image reconstruction. Nature Methods, 15(5):309, 2018.
- [45] A. M. Stuart. Inverse problems: A Bayesian perspective. Acta Numer., 19:451–559, 2010.
- [46] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *Proceedings of the Int. Conf. on Learning Representations*, 2014.
- [47] E. van den Berg and M. P. Friedlander. Probing the pareto frontier for basis pursuit solutions. SIAM Journal on Scientific Computing, 31(2):890–912, 2008.
- [48] E. Weinan and B. Yu. The Deep Ritz Method: A Deep Learning-Based Numerical Algorithm for Solving Variational Problems. *Commun. Math. Stat.*, 6(1):1–14, 2018.
- [49] D. Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. arXiv:1802.03620, 2018.
- [50] B. Zhu, J. Z. Liu, S. F. Cauley, B. R. Rosen, and M. S. Rosen. Image reconstruction by domain-transform manifold learning. *Nature*, 555(7697):487, 03 2018.