

Title

Optional Subtitle

Andras Filip Plaian

Master's Thesis, Autumn 2021



This master's thesis is submitted under the master's program *Computational Science*, with program option *Applied Mathematics and Risk Analysis*, at the Department of Mathematics, University of Oslo. The scope of the thesis is 60 credits.

The front page depicts a section of the root system of the exceptional Lie group E_8 , projected into the plane. Lie groups were invented by the Norwegian mathematician Sophus Lie (1842–1899) to express symmetries in differential equations and today they play a central role in various parts of mathematics.

Abstract

Add new section about results in

Acknowledgements

Rewrite this.

Contents

Abstract	i
Acknowledgements	iii
Contents	v
List of Figures	vii
1 Introduction	1
2 Concepts From Numerical Analysis	3
2.1 Numerical Stability	3
2.2 Numerical Accuracy	3
2.3 The Importance Of Numerical Stability	4
3 Compressed Sensing	5
3.1 Main Assumptions: Sparsity and Compressibility	5
3.2 Basic Requirements of A	6
3.3 Finding Suitable Reconstruction Algorithms	8
3.4 Basis Pursuit	12
3.5 Matrix Analysis: Coherence	17
4 Neural Networks And Deep Learning	19
4.1 Supervised Machine Learning	19
4.2 Design	20
4.3 Training	21
4.4 The Universal Approximation Theorem	22
4.5 The Success Of Deep Learning	23
4.6 Instabilities In Deep Learning	24
4.7 Consequences of the false structure phenomenon	25
4.8 Reasons To Go Beyond Compressed Sensing	26
4.9 Deep Learning for inverse problems	26
5 Instabilities In Deep Learning For Inverse Problems	29
5.1 Universal Instabilities In Inverse Reconstruction Methods . . .	29
5.2 False Positives	30
5.3 False Negatives	32

Contents

5.4 Stability Through Kernel Awareness	33
Appendices	37
Bibliography	39

List of Figures

4.1	Extremely simplified NN with a single hidden layer.	21
4.2	Different choices of activations functions.	21
4.3	3×3 images with either horizontal or vertical stripes.	25
5.1	False Positive, images from [11].	31
5.2	False Negatives, images from [11].	33
5.3	CS stable recovery vs unstable DL recovery, figures from [11]. . . .	34
5.4	CS stable recovery vs unstable DL recovery, figures from [11]. . . .	35

CHAPTER 1

Introduction

Write introduction here.

Outline

The rest of the text is organised as follows:

CHAPTER 2

Concepts From Numerical Analysis

We begin this chapter by introducing a few concepts from numerical analysis, that will be used throughout the thesis in different contexts.

2.1 Numerical Stability

An important part of numerical analysis is to figure out if a given algorithm is *numerically* stable or not. An algorithm is called numerically stable if an error, whatever its cause, does not grow to be much larger during the calculation [1]. In order to estimate the bounds to these errors, there exist the following notion of a Lipschitz constant:

Definition 2.1 (Lipschitz constants). Let (X, d_X) and (Y, d_Y) be metric spaces, and let the function $f : X \rightarrow Y$ be Lipschitz continuous such that there exist a number $K \in \mathbb{R}$ that satisfy

$$d_Y(f(x), f(y)) \leq K d_X(x, y) \text{ for all } x, y \in X. \quad (2.1)$$

The smallest possible K that satisfy (2.1), is the said to be the Lipschitz constant of f .

As we can see from the definition above, the Lipschitz constant gives an upper bound for how much an error in the input of a function may change the error in the output of the function. Depending on the value of K , we have three ways the function f can change the input error:

- $K < 1$, the error reduces, such a function is called a *contraction*.
- $K = 1$, the error stays the same.
- $K > 1$, the error increases.

2.2 Numerical Accuracy

Accuracy refers to the absolute or relative error of an approximate quantity [1]. Suppose we have a numerical method for estimating the true solution x , then the absolute error ϵ caused by our estimate x^* , is given by:

$$x^* = x + \epsilon \quad (2.2)$$

2. Concepts From Numerical Analysis

While the relative error δ is given by:

$$x^* = x(1 + \delta) \tag{2.3}$$

2.3 The Importance Of Numerical Stability

During the exploding interest of Machine Learning in the last decades, the community have suffered from some growing pains. As there have been a huge push to publish new methods, in many cases being the first is more important than making sure all the details are correct [2].

Sometimes we need to be reminded that with the development of new technology, we must make sure that the methods we develop are secure, particularly for applications where the lack of stable and accurate methods may result into unfortunate events. Some of these unfortunate events, that resulted from unstable numerical methods, include The Sleipner A offshore platform accident in 1991, causing an estimated loss of 1.8 billion NOK [3], The Patriot Missile Failure in 1998, killing 28 and injuring 98 people [4], and The Explosion of the Ariane 5 costing 500 million USD [5].

CHAPTER 3

Compressed Sensing

In an underdetermined system of linear equations there are fewer equations than unknowns. In mathematical terms this can be stated as the matrix equation $A\mathbf{x} = \mathbf{y}$ where $A \in \mathbb{C}^{m \times N}$, $\mathbf{x} \in \mathbb{C}^N$, $\mathbf{y} \in \mathbb{C}^m$ and $m < N$. From classical linear algebra it follows that this equation is unsolvable in the general case, however, given certain conditions, it is possible to find an exact or an estimate of a solution.

The research area associated with these assumptions is called Compressed Sensing (CS). The goal of this chapter is to introduce the reader to the needed concepts in order to be able to determine the required conditions for having a stable and robust reconstruction method, for finding solutions of the *inverse problem* $\mathbf{x} = A^{-1}\mathbf{y}$ by exploiting the theory of CS. Most of the material presented in this chapter is based upon the book of Foucart and Rauhut, which is a comprehensive introduction to Compressed Sensing, thus where details are lacking, the reader is referred to [6]. See also the work of Adcock and Hansen [7], for an in-depth treatment.

After this chapter, roughly speaking, we should be able to have some answers to the following questions:

- What properties does A need to possess in order for \mathbf{x} to be reconstructable?
- Does feasible algorithms even exist for solving the inverse problem?

3.1 Main Assumptions: Sparsity and Compressibility

The main assumption we have about the vector \mathbf{x} , which we are trying to recover, is that it is *sparse*. Informally, a vector is sparse if most of its components are zero. Although, a change of basis might be needed in order for \mathbf{x} to be sparse in the problem of interest. For instance in image reconstruction, the natural images are sparse in the wavelet domain. Interestingly, natural scenes sparsely activate neurons in the primary visual cortex in humans [8].

Before we define sparsity formally, we recall the definition of the *support* of a given vector:

3. Compressed Sensing

Definition 3.1. The support of a vector $\mathbf{x} \in \mathbb{C}^N$ is the index set of its non-zero entries, that is:

$$\text{supp}(\mathbf{x}) := \{i \in [N] : x_i \neq 0\}.$$

In Compressive Sensing it is customary to abuse the l_0 notation to denote the cardinality of the support set, i.e. the number of non-zero entries of a vector. The $\|\cdot\|_0$ fails to be a norm since it does not satisfy the scaling property of norms. We can now define *s-sparse* vectors.

Definition 3.2. The vector $\mathbf{x} \in \mathbb{C}^N$ is called *s-sparse* if it has no more than s non-zero entries, that is if: $\|\mathbf{x}\|_0 \leq s$.

The notion of sparsity is an ideal one, meaning that in the real world it may very well be that the vector we are trying to recover is only close to being sparse. In order for CS to handle more problems, we may also be interested in the following notion of *compressibility*.

Definition 3.3. For $p > 0$, the measure of a vectors compressibility is given by the l_p -error of best s -term approximation to $\mathbf{x} \in \mathbb{C}^N$ defined by

$$\sigma(\mathbf{x})_p := \inf\{\|\mathbf{x} - \mathbf{z}\|_p, \mathbf{z} \in \mathbb{C}^N \text{ is } s\text{-sparse}\}$$

Informally, a vector is compressible if the l_p -error decays quickly in s . For instance in image processing, if we can describe most of the variance in image data with a few components of a vector.

3.2 Basic Requirements of A

What properties does A need to have in order to have the possibility to reconstruct every s -sparse vector *regardless* of the reconstruction method? For instance, what is the minimal number of rows that A needs to have such that it would be possible to find every s -sparse vector? We have the following result:

Theorem 3.4. For a given matrix $A \in \mathbb{C}^{m \times N}$, the following statements are equivalent:

- a) Every s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ is the unique s -sparse solution of $A\mathbf{z} = A\mathbf{x}$, that is, if $A\mathbf{x} = A\mathbf{z}$ and both \mathbf{z} and \mathbf{x} are s -sparse, then $\mathbf{z} = \mathbf{x}$.
- b) The null space $\ker(A)$ does not contain any $2s$ -sparse vector other than the zero vector, that is, $\ker(A) \cap \{\mathbf{z} \in \mathbb{C}^N : \|\mathbf{z}\|_0 \leq 2s\} = \{\mathbf{0}\}$.
- c) For every $S \subset [N]$ with $\text{card}(S) \leq 2s$, the submatrix A_S is injective.
- d) Every set of $2s$ columns of A is linearly independent.

Proof. b) \implies a) Let \mathbf{x}, \mathbf{z} be s -sparse vectors, not necessarily unique, such that $A\mathbf{x} = A\mathbf{z}$. Then $\mathbf{x} - \mathbf{z}$ is at most $2s$ -sparse and using linearity of A we get $A(\mathbf{x} - \mathbf{z}) = \mathbf{0}$. If b) holds, then the kernel only contains the zero vector of $2s$ -sparse vectors, then the only way $A(\mathbf{x} - \mathbf{z}) = \mathbf{0}$ is true, is if $\mathbf{x} = \mathbf{z}$, thus \mathbf{x} is the unique s -sparse vector.

a) \implies b) Assume that \mathbf{x} is the unique s -sparse vector such that $A\mathbf{x} = A\mathbf{z}$. Let $\mathbf{v} \in \ker(A)$ be $2s$ -sparse. Since \mathbf{v} is $2s$ -sparse, there exist a way to construct \mathbf{v} from $\mathbf{x} - \mathbf{z}$ where \mathbf{x}, \mathbf{z} are both s -sparse and such that they have no components in common. Assuming a) holds, then we have $\mathbf{x} = \mathbf{z} = \mathbf{v} = \mathbf{0}$.

b) \iff c) \iff d) For a $2s$ -sparse vector \mathbf{v} with $S = \text{supp } \mathbf{v}$, if $\mathbf{a}_1, \dots, \mathbf{a}_S$ are the column vectors of the submatrix A_S made from A , we get $A\mathbf{v} = A_S\mathbf{v}_S$. Noting that $S = \text{supp } \mathbf{v}$ ranges through all possible subsets of $[N]$ with $\text{card}(S) \leq 2s$ ranging through all the possible $2s$ -sparse \mathbf{v} vectors. Then in this view, linear independence d), injectivity c) and that the kernel is trivial b) are all equivalent statements about A . ■

With a little knowledge about how to construct invertible matrices, we can derive the following result from the above theorem:

Corollary 3.5. For any integer $N \geq 2s$, there exists a matrix $A \in \mathbb{C}^{m \times N}$ with $m = 2s$ rows such that every s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ can be recovered from the vector $\mathbf{y} = A\mathbf{x} \in \mathbb{C}^m$.

Proof. Let t_1, t_2, \dots, t_N be strictly positive numbers and consider the matrix $A \in \mathbb{C}^{m \times N}$ with $m = 2s$ defined by

$$A = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ t_1 & t_2 & \cdots & t_N \\ \vdots & \vdots & \cdots & \vdots \\ t_1^{2s-1} & t_2^{2s-1} & \cdots & t_N^{2s-1} \end{bmatrix}$$

3. Compressed Sensing

Now define a square submatrix $A_S \in \mathbb{C}^{m \times m}$ indexed by $S = \{i_1 < \dots < i_m\}$ as

$$A_S = \begin{bmatrix} 1 & 1 & \dots & 1 \\ t_{i_1} & t_{i_2} & \dots & t_{i_m} \\ \vdots & \vdots & \ddots & \vdots \\ t_{i_1}^{2s-1} & t_{i_2}^{2s-1} & \dots & t_{i_m}^{2s-1} \end{bmatrix}$$

The given submatrix A_S can be recognized as the transposed *Vandermonde matrix*. A result about Vandermonde matrices is that its determinant, in view of A_S , can be expressed as $\det(A_S) = \prod_{k < l \leq m} (t_{i_l} - t_{i_k})$. Since the t 's were defined to be strictly positive it follows that $\det(A_S) > 0$, and hence A_S is invertible. The Invertible Matrix Theorem then tells us that A_S is injective which means that condition *iii*) in Theorem 2.4 holds and then we have a unique s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ such that $A\mathbf{z} = A\mathbf{x}$. ■

3.3 Finding Suitable Reconstruction Algorithms

Recall that the compressed sensing problem is to find an s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ given $\mathbf{y} = A\mathbf{x}$. Abusing the customary l_0 notation for the number of non-zero entries, we can reformulate this problem as an optimization problem, that is:

$$\text{minimize } \|\mathbf{z}\|_0 \quad \text{subject to } A\mathbf{z} = \mathbf{y}. \quad (P_0)$$

A naive approach to solving this could be to look at it as a combinatorial optimization problem and use brute force to calculate every square system $A_S^* A_S \mathbf{x} = A_S^* \mathbf{y}$, for $\mathbf{x} \in \mathbb{C}^S$ where S runs through all possible subsets of $[N]$ with size s . This might very well work on small sizes of N , but the total amount of subsets the algorithm has to go through is determined by the formula $\binom{N}{s}$. A quick illustration shows that with $N = 1000, s = 10$ we have $\binom{1000}{10} \geq \left(\frac{1000}{10}\right)^{10} = 10^{20}$ linear systems of size 10×10 to solve. Assuming one could solve each system in 10^{-10} seconds, it would still take 10^{10} seconds, ie. several human lifespans, thus the brute force approach is completely unpractical for sufficiently large N .

P_0 is NP-hard

It's not only the brute force approach that is not bounded by a polynomial expression, it might be that there doesn't exist an approach bounded by a polynomial expression at all. In fact, we have the following result:

Theorem 3.6. *The l_0 -minimization problem for an arbitrary matrix $A \in \mathbb{C}^{m \times N}$ and a vector $\mathbf{y} \in \mathbb{C}^m$ is NP-hard.*

Before proving the result, let us recall some of the terminology:

- The class of P problems consists of all decision problems for which there exists a polynomial-time algorithm, ie. an input size bounded by a polynomial expression, for finding a solution.

3.3. Finding Suitable Reconstruction Algorithms

- The class of NP problems, not to be confused with NP-hard, consists of all decision problems for which there exists a polynomial-time algorithm *certifying* a solution.
- The class of NP-hard problems consists of all problems for which *solving* algorithm could be transformed in polynomial time into a *solving* algorithm for any NP-problem.

This means that if we can show that a given problem solves an other problem belonging a certain class of problems, then the given problems also belongs to that same class.

Proof. If we can show that a known NP-hard problem can be reduced in polynomial time to the l_0 -minimization problem then l_0 -minimization itself is NP-hard.

Let The Exact Cover by 3-sets problem be our known NP-hard problem. The Exact Cover by 3-sets problem says that given a collection $\{\mathcal{C}_i, i \in [N]\}$ of 3-element subsets of $[m]$, does there exist an exact partition or cover of the set $\{1, 2, \dots, m\} = [m]$?

Let $\{\mathcal{C}_i, i \in [N]\}$ be the collection of 3-element subsets of $[m]$. Define vectors $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{C}^m$ by

$$a_{ij} = \begin{cases} 1 & \text{if } j \in \mathcal{C}_i, \\ 0 & \text{if } j \notin \mathcal{C}_i. \end{cases}$$

where $j \in [m]$.

Define a matrix $A \in \mathbb{C}^{m \times N}$ and a vector $\mathbf{y} \in \mathbb{C}^m$ by

$$A = [\mathbf{a}_1 \cdots \mathbf{a}_N] \text{ and } \mathbf{y} = [1, \dots, 1]^T.$$

Since $N \leq \binom{m}{3}$, this construction can be done in polynomial time since $\binom{m}{3} = \frac{m(m-1)(m-2)}{3!}$ which is a 3rd degree polynomial in m .

If a vector $\mathbf{z} \in \mathbb{C}^N$ satisfies $A\mathbf{z} = \mathbf{y}$, then all the m components of $A\mathbf{z}$ are nonzero and $\|A\mathbf{z}\|_0 = m$. Since each vector \mathbf{a}_i has exactly 3 nonzero components, the vector $A\mathbf{z} = \sum_{j=1}^N z_j \mathbf{a}_j$ has at most 3 $\|\mathbf{z}\|_0$ nonzero components, $\|A\mathbf{z}\|_0 \leq 3\|\mathbf{z}\|_0$ and consequently $\|\mathbf{z}\|_0 \geq m/3$. Now consider the l_0 -minimization and let $\mathbf{x} \in \mathbb{C}^N$ be the output. We separate two cases:

1. If $\|\mathbf{x}\|_0 = m/3$, then the collection $\{\mathcal{C}_j, j \in \text{supp}(\mathbf{x})\}$ forms an exact cover of $[m]$, for otherwise the m components of $A\mathbf{x} = \sum_{j=1}^N x_j \mathbf{a}_j$ would not all be nonzero.
2. If $\|\mathbf{x}\|_0 > m/3$, then no exact cover $\{\mathcal{C}_j, j \in J\}$ can exist, for otherwise the vector $\mathbf{z} \in \mathbb{C}^N$, defined by $z_j = 1$ if $j \in J$ and $z_j = 0$ if $j \notin J$, would satisfy $A\mathbf{z} = \mathbf{y}$ and $\|\mathbf{z}\|_0 = m/3$, contradicting the l_0 -minimality of \mathbf{x} .

This shows that solving l_0 -minimization problem enables one to solve the Exact Cover by 3-sets problem and consequently l_0 -minimization problem is NP-hard. ■

3. Compressed Sensing

P_q where $0 < q < 1$ also NP-hard

If P_0 optimization is intractable, would P_q optimization be a better approach? Not according to the following proposition:

Proposition 3.7. l_q -minimization for $0 < q < 1$ is NP-hard.

Proof. The proof goes along the same lines as proving the NP-hardness of l_0 -minimization, that is if l_q -minimization can help to verify that a given NP-hard problem has a solution, that is if the minimum of $\|\mathbf{w}\|_q$ subject to $A\mathbf{w} = \mathbf{y}$, equals n , then l_q -minimization must necessarily also belong to the same class of NP-hard or NP-complete problems.

The *partition problem* consists, given integers a_1, a_2, \dots, a_n , in deciding whether there exists two sets $I, J \subset [n]$ such that $I \cap J = \emptyset$, $I \cup J = [n]$ and $\sum_{i \in I} a_i = \sum_{j \in J} a_j$ for $i, j \in I, J$ respectively.

Define

$$A = \begin{bmatrix} a_1 & a_2 & \cdots & a_n & -a_1 & -a_2 & \cdots & -a_n \\ 1 & 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & 0 & 1 & \cdots & 0 \\ \vdots & & \ddots & 0 & \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & 1 & 0 & \cdots & 0 & 1 \end{bmatrix} \text{ and } \mathbf{y} = [0, 1, 1, \dots, 1]^T.$$

Let \mathbf{x} and \mathbf{z} be the first and second half of the input vector \mathbf{w} to A respectively. That $A\mathbf{w} = \mathbf{y}$ implies that $\sum_{i=1}^n a_i x_i = \sum_{i=1}^n a_i z_i$ and $x_i + z_i = 1$ for all i . We seek to minimize $\sum_{i=1}^n (x_i^q + z_i^q)$ under these constraints.

The minimum value for the objective function when only the constraints $x_i + z_i$ are considered is n , and this occurs if and only if only 0's and 1's are involved. If the minimum $\|\mathbf{w}\|_q$ subject to $A\mathbf{w} = \mathbf{y}$ is n , then it will also be n when only the $x_i + z_i = 1$ constraints are considered. This means that either x_i or z_i is equal to 1, for all i . Let I be the set of i 's where $x_i=1$, and J be the set of i 's where $z_i=1$. We then have $I \cap J = \emptyset$ and $I \cup J = [n]$, ie. a partition on the form we seek.

Defining x_i as 1 when $i \in I$ and defining z similarly, we obtain a vector \mathbf{w} where all the constraints are fulfilled and where the minimum n is obtained. ■

P_q where $q > 1$

If both l_0 and l_q -minimization for $0 < q < 1$ is NP-hard, how about l_q -minimization for $q > 1$? Unfortunately, this optimization problem fails to recover even 1-sparse vectors. We have the following proposition:

Proposition 3.8. Let $q > 1$ and let A be a $m \times N$ matrix with $m < N$. Then there exists a 1-sparse vector which is not a minimizer of P_q .

Proof. Since $m < N$, we have from linear algebra that the kernel of A is non-trivial. Hence there must exist a $\mathbf{v} \neq \mathbf{0} \in \ker(A)$ such that $A\mathbf{v} = \mathbf{0}$.

3.3. Finding Suitable Reconstruction Algorithms

Assume, for the sake of contradiction, that all standard basis vectors \mathbf{e}_j are minimizers. Choose a j so that $v_j \neq 0$. Then

$$\|\mathbf{e}_j + t\mathbf{v}\|_q^q = |1 + tv_j|^q + \sum_{k \neq j} |tv_k|^q = |1 + tv_j|^q + |t|^q \sum_{k \neq j} |v_k|^q$$

Define two functions such that

$$g_+(t) = (1 + tv_j)^q + t^q \sum_{k \neq j} |v_k|^q$$

$$g_-(t) = (1 + tv_j)^q + (-t)^q \sum_{k \neq j} |v_k|^q$$

For $|t| < 1/v_j$, $\|\mathbf{e}_j + t\mathbf{v}\|_q^q$ coincides with g_+ for $t \geq 0$, and g_- for $t \leq 0$. Taking the derivatives of g_+ and g_- we get

$$g'_+(t) = qv_j(1 + tv_j)^{q-1} + qt^{q-1} \sum_{k \neq j} |v_k|^q$$

$$g'_-(t) = qv_j(1 + tv_j)^{q-1} - q(-t)^{q-1} \sum_{k \neq j} |v_k|^q$$

Taking the limit of the two derivatives and then for $q > 1$ we get

$$\lim_{t \rightarrow 0^+} g'_+(t) = \lim_{t \rightarrow 0^-} g'_-(t) = g'_-(t) = qv_j$$

This means that near 0, $\|\mathbf{e}_j + t\mathbf{v}\|_q^q$ has derivative arbitrarily near $qv_j \neq 0$. Since a linear function with derivative qv_j has no minimum near 0, then with $t = 0$, \mathbf{e}_j can't be a minimizer of P_q , which consequently contradicts our assumption that all standard basis vectors \mathbf{e}_j are minimizers of P_q . ■

The special case P_q where $q = 1$

Recall that in an optimization problem we try to minimize or maximize an objective function subject to constraint functions. If all the constraint functions are also convex functions, then the optimization problem becomes a convex optimization problem. Thus both l_0 and l_q -minimization for $0 < q < 1$ are non-convex optimization problems. If all the constraint functions are linear functions, then the following convex optimization problem

$$\text{minimize } \|\mathbf{z}\|_1 \quad \text{subject to } A\mathbf{z} = \mathbf{y}. \quad (P_1)$$

is also a linear program. Efficient algorithms exist for solving LP's. For instance, the Simplex method, behaves like a polynomial-time algorithm for solving real-life LP problems. l_1 -minimization is also known *basis pursuit* and during the next sections we shall see that given certain assumptions on A we can actually solve P_1 and that the recovered solutions will be sparse.

3. Compressed Sensing

3.4 Basis Pursuit

In order to show how BP can identify the solution to (P_0) , we need to introduce some definitions and theorems, an important property which our matrix must have, is the so-called Null Space Property of a matrix A :

Definition 3.9. A matrix $A \in \mathbb{C}^{m \times N}$ is said to satisfy the *Null Space Property (NSP)* relative to a set $S \subset \{1, 2, \dots, N\}$ if

$$\min \|\mathbf{v}_S\|_1 < \|\mathbf{v}_{\bar{S}}\|_1 \text{ for all } \mathbf{v} \in \ker A \setminus \{\mathbf{0}\}$$

A is said to satisfy the *Null Space Property of order s* if it satisfies the null space property relative to any set $S \subset \{1, 2, \dots, N\}$ with $|S| \leq s$.

The following theorem shows that the NSP of a matrix is a sufficient condition in order to solve (P_0) .

Theorem 3.10. *Given a matrix $A \in \mathbb{C}^{m \times N}$, every vector $\mathbf{x} \in \mathbb{C}^N$ supported on a set S is the unique solution (P_1) with $A\mathbf{x} = \mathbf{y}$ if and only if A satisfies the NSP relative to S .*

Furthermore, if the set S varies, then every s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ is the unique solution to (P_1) with $A\mathbf{x} = \mathbf{y}$ if and only if A satisfies the NSP of order s .

Proof. Let S be a fixed index set, and assume that every vector $\mathbf{x} \in \mathbb{C}^N$ supported on this set, is the unique minimizer of (P_1) . From the assumption it follows that for $\mathbf{v} \in \ker A \setminus \{\mathbf{0}\}$, the vector \mathbf{v}_S is the unique minimizer of (P_1) . Since $A(\mathbf{v}_S + \mathbf{v}_{\bar{S}}) = \mathbf{0}$ and $-\mathbf{v}_S \neq \mathbf{v}_{\bar{S}}$, from the minimality assumption we must have that $\|-\mathbf{v}_S\|_1 < \|\mathbf{v}_{\bar{S}}\|_1$. This established the NSP relative to S .

Conversely, assume that NSP relative to S holds. Let $\mathbf{x} \in \mathbb{C}^N$ be supported on S and a vector $\mathbf{z} \in \mathbb{C}^N$, $\mathbf{z} \neq \mathbf{x}$, such that $A\mathbf{z} = A\mathbf{x}$. Following the rules for norms and taking complements for the support of a set, we obtain

$$\|\mathbf{x}\|_1 \leq \|\mathbf{x} - \mathbf{z}_S\|_1 + \|\mathbf{z}_S\|_1 = \|\mathbf{v}_S\|_1 + \|\mathbf{z}_S\|_1 < \|\mathbf{v}_{\bar{S}}\|_1 + \|\mathbf{z}_S\|_1 = \|-\mathbf{z}_{\bar{S}}\|_1 + \|\mathbf{z}_S\|_1 = \|\mathbf{z}\|_1.$$

Which shows that \mathbf{x} obtains the unique minimum.

To prove the second part of the theorem, let S vary and assume that every s -sparse vector \mathbf{x} is found by solving (P_1) subject to $A\mathbf{x} = \mathbf{y}$. Let \mathbf{z} be the solution to P_0 subject to $A\mathbf{x} = \mathbf{y}$ then $\|\mathbf{z}\|_0 \leq \|\mathbf{x}\|_0$ so that also \mathbf{z} is s -sparse. But since every s -sparse vector is the unique minimizer of (P_1) , we have that $\mathbf{x} = \mathbf{z}$ and the result follows. ■

Minimum Number Of Rows For Basis Pursuit

From the results above it is clear that if a matrix possesses the NSP of order s , the BP will solve (P_1) . Next we will introduce a theorem that can identify when A has the NSP of order s .

Theorem 3.11. *Given a matrix $A \in \mathbb{C}^{m \times N}$, then every set of $2s$ columns of A is linearly independent if and only if A satisfies the NSP of order s .*

Proof. Assume that every $2s$ columns of A is linearly independent, then from The Invertible Matrix Theorem, we have that the kernel of A does not contain any other $2s$ -sparse vector other than $\mathbf{0}$. Now let \mathbf{x} , and \mathbf{z} be s -sparse with $A\mathbf{z} = A\mathbf{x}$. Then $A(\mathbf{x} - \mathbf{z}) = \mathbf{0}$ and $\mathbf{x} - \mathbf{z}$ is $2s$ -sparse, but since $\ker A \setminus \{\mathbf{0}\}$ is empty, we must have $\mathbf{x} = \mathbf{z}$, but this implies that the NSP of order s holds. Conversely, assume that the kernel of A does not contain any other $2s$ -sparse vector other than $\mathbf{0}$, then for any set S with $\text{card}(S) = 2s$, $A_S \mathbf{x} = \mathbf{0}$ only has the trivial solution and thereby $2s$ linearly independent columns. ■

From these results we can derive that :

Corollary 3.12. *In order for Basis Pursuit to solve (P_0) , the matrix $A \in \mathbb{C}^{m \times N}$ has to satisfy:*

$$m \geq 2s$$

Proof. Assume that it is possible to uniquely recover any s -sparse vector \mathbf{x} from $\mathbf{y} = A\mathbf{x}$. Then, by Theorem 2.4, statement a) holds, and consequently so does d), that is, that every set of $2s$ columns of A is linearly independent. Thus $\text{rank}(A) \geq 2s$, but we also have that the rank of a matrix can not be greater than the number of rows m , so we must have $\text{rank}(A) \leq m$. Combining and arranging these two inequalities, we get that

$$2s \leq \text{rank}(A) \leq m$$

which implies, $m \geq 2s$, the inequality in the corollary. ■

Stability

Basis Pursuit features another important property, namely *stability*. This property tackles sparsity defects, that is, to recover a vector $\mathbf{x} \in \mathbb{C}^N$ whose error to the true underlying vector is controlled by its distance to an s -sparse vector.

Definition 3.13. A matrix $A \in \mathbb{C}^{m \times N}$ is said to satisfy the *stable null space property* with $0 < \rho < 1$ relative to a set $S \subset [N]$ if

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1 \text{ for all } \mathbf{v} \in \ker A$$

Furthermore, it satisfies the stable null space of order s if it satisfies the stable null space property relative to any set $S \subset [N]$ with $\text{card}(S) \leq s$.

If a matrix satisfies the stable null space property of order s , then the l_1 -error to the true underlying vector \mathbf{x} is given by

$$\|\mathbf{x} - \mathbf{x}'\|_1 \leq \frac{2(1 + \rho)}{1 - \rho} \sigma_s(\mathbf{x})_1 \quad (3.1)$$

3. Compressed Sensing

The next result will give us a condition for when, the inequality (2.1) and the stable null space property, will hold.

Theorem 3.14. *The matrix $A \in \mathbb{C}^{m \times N}$ satisfies the stable null space property with constant $0 < \rho < 1$ relative to S if and only if*

$$\|\mathbf{z} - \mathbf{x}\|_1 \leq \frac{1 + \rho}{1 - \rho} (\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1) \quad (3.2)$$

for all vectors $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$ with $A\mathbf{z} = A\mathbf{x}$.

Before proving Theorem 2.14, we show that (2.2) implies (2.1). Let S be the set of the s largest non-negative components of \mathbf{x} , then $\|\mathbf{x}_{\bar{S}}\|_1 = \sigma_s(\mathbf{x})_1$. If \mathbf{x}' is a minimizer of (P_1) , then $\|\mathbf{x}'\|_1 \leq \|\mathbf{x}\|_1$ and $A\mathbf{x}' = A\mathbf{x}$. The right-hand side of the inequality (2.2) with $\mathbf{z} = \mathbf{x}'$ will then be equal to the right hand side of inequality (2.1).

Proof. To prove Theorem 2.14, assume that the matrix A satisfies inequality (2.2) for all vectors $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$ with $A\mathbf{z} = A\mathbf{x}$. Given a vector $\mathbf{v} \in \ker(A)$, since $A\mathbf{v}_{\bar{S}} = A(-\mathbf{v}_S)$, we can apply (2.2) with $\mathbf{x} = -\mathbf{v}_S$ and $\mathbf{z} = \mathbf{v}_{\bar{S}}$. This gives us that

$$\|\mathbf{v}\|_1 \leq \frac{1 + \rho}{1 - \rho} (\|\mathbf{v}_S\|_1 - \|\mathbf{v}_{\bar{S}}\|_1)$$

Rearranging and cancelling equal terms, this can be written as

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1$$

which we recognize as the stable NSP with constant $0 < \rho < 1$.

Conversely, assume that the matrix A satisfies the stable NSP with constant $0 < \rho < 1$ relative to S . Then for $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$ with $A\mathbf{z} = A\mathbf{x}$, we get that $\mathbf{v} = \mathbf{z} - \mathbf{x} \in \ker(A)$ yields

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1. \quad (3.3)$$

Since $\|\mathbf{v}_{\bar{S}}\|_1 = \|(\mathbf{z} - \mathbf{x})_{\bar{S}}\|_1$, we can rewrite (2.3) into

$$\|\mathbf{v}_{\bar{S}}\|_1 \leq \|\mathbf{z}_{\bar{S}}\|_1 - \|\mathbf{x}_{\bar{S}}\|_1 + \rho \|\mathbf{v}_{\bar{S}}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1.$$

Since $\rho < 1$, this can be rewritten as

$$\|\mathbf{v}_{\bar{S}}\|_1 \leq \frac{1}{1 - \rho} (\|\mathbf{z}_{\bar{S}}\|_1 - \|\mathbf{x}_{\bar{S}}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1).$$

Using (2.3) once again, we chain the inequalities to get,

$$\|\mathbf{v}\|_1 = \|\mathbf{v}_{\bar{S}}\|_1 + \|\mathbf{v}_S\|_1 \leq (1 + \rho) \|\mathbf{v}_{\bar{S}}\|_1 \leq \frac{1 + \rho}{1 - \rho} (\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1)$$

which is (2.2), the desired inequality. ■

Robustness

Basis Pursuit features another important property, namely *robustness*. This property tackles perturbations in the input.

Definition 3.15. Given $q \geq 1$, the matrix $A \in \mathbb{C}^{m \times N}$ is said to satisfy the l_q -robust null space property (rNSP), with constants $0 < \rho < 1$ and $\tau > 0$, relative to a set $S \subset [N]$ if

$$\|\mathbf{v}_S\|_q < \frac{\rho}{s^{1-1/q}} \|\mathbf{v}_{\bar{S}}\|_1 + \tau \|A\mathbf{v}\| \text{ for all } \mathbf{v} \in \mathbb{C}^N.$$

Furthermore, it satisfies the l_q -rNSP of order s if it satisfies the l_q -rNSP relative to any set $S \subset [N]$ with $\text{card}(S) \leq s$.

In order to introduce a more general result using l_q -rNSP, we need an intermediate result that holds for the ordinary rNSP, that is setting $q = 1$.

Theorem 3.16. The matrix $A \in \mathbb{C}^{m \times N}$ satisfies the robust null space property with constants $0 < \rho < 1$ and $\tau > 0$ relative to S if and only if

$$\|\mathbf{z} - \mathbf{x}\|_1 \leq \frac{1+\rho}{1-\rho} (\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1) + \frac{2\tau}{1-\rho} \|A(\mathbf{z} - \mathbf{x})\| \quad (3.4)$$

for all vectors $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$.

Proof. Very similar to the proof of Theorem 2.14, will therefore use results from there when appropriate.

Assume A satisfies (2.4) for all vectors $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$. Let $\mathbf{v} \in \mathbb{C}^N$ be so that $\mathbf{v} = \mathbf{z} - \mathbf{x}$ with $\mathbf{z} = \mathbf{v}_{\bar{S}}$ and $\mathbf{x} = -\mathbf{v}_S$ which gives

$$\|\mathbf{v}\|_1 \leq \frac{1+\rho}{1-\rho} (\|\mathbf{v}_{\bar{S}}\|_1 - \|\mathbf{v}_S\|_1) + \frac{2\tau}{1-\rho} \|A\mathbf{v}\|.$$

Multiplying by $(1-\rho)$ on both sides and expanding $\|\mathbf{v}\|_1$ on the left side gives

$$(1-\rho)(\|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1) \leq (1+\rho)(\|\mathbf{v}_{\bar{S}}\|_1 - \|\mathbf{v}_S\|_1) + 2\tau \|A\mathbf{v}\|.$$

Expanding, cancelling equal terms and dividing by 2 on both sides results in

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1 + \tau \|A\mathbf{v}\|$$

which is the rNSP with constants $0 < \rho < 1$ and $\tau > 0$.

For the converse, assume A satisfies rNSP with constants $0 < \rho < 1$ and $\tau > 0$ relative to S . For $\mathbf{z}, \mathbf{x} \in \mathbb{C}^N$ with $\mathbf{v} = \mathbf{z} - \mathbf{x}$, the rNSP gives

$$\|\mathbf{v}_S\|_1 \leq \rho \|\mathbf{v}_{\bar{S}}\|_1 + \tau \|A\mathbf{v}\|,$$

$$\|\mathbf{v}_{\bar{S}}\|_1 \leq \|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + \|\mathbf{v}_S\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1$$

3. Compressed Sensing

Combining these two inequalities and using the rNSP again we arrive at

$$\begin{aligned}\|\mathbf{v}\|_1 &= \|\mathbf{v}_S\|_1 + \|\mathbf{v}_{\bar{S}}\|_1 \leq (1 + \rho)\|\mathbf{v}_S\|_1 + \tau\|A\mathbf{v}\| \\ &\leq \frac{1 + \rho}{1 - \rho}(\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\|\mathbf{x}_{\bar{S}}\|_1) + \frac{2\tau}{1 - \rho}\|A\mathbf{v}\|,\end{aligned}$$

which is the desired inequality. ■

Theorem 3.17. *Given $1 \leq p \leq q$, suppose that the matrix $A \in \mathbb{C}^{m \times N}$ satisfies the l_q -robust null space property of order s with constants $0 < \rho < 1$ and $\tau > 0$. Then, for any $\mathbf{x}, \mathbf{z} \in \mathbb{C}^N$,*

$$\|\mathbf{x} - \mathbf{z}\|_p \leq \frac{C}{s^{1-1/p}}(\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\sigma_s(\mathbf{x})_1) + Ds^{1/p-1/q}\|A(\mathbf{z} - \mathbf{x})\|,$$

where $C := (1 + \rho)^2/(1 - \rho)$ and $D := (3 + \rho)\tau/(1 - \rho)$.

In order to prove the above result, we need the following little lemma:

Lemma 3.18. *For any $q > p > 0$ and any $\mathbf{x} \in \mathbb{C}^N$, the inequality*

$$\sigma_s(\mathbf{x})_q \leq \frac{c_{p,q}}{s^{1/p-1/q}}\|\mathbf{x}\|_p$$

holds with

$$c_{p,q} = \left[\left(\frac{p}{q} \right)^{p/q} \left(1 - \frac{p}{q} \right)^{1-p/q} \right]^{1/p} \leq 1.$$

Particularly for our proof, with the choice $p = 1$ and $q = p$ gives

$$\sigma_s(\mathbf{x})_p \leq \frac{1}{s^{1-1/p}}\|\mathbf{x}\|_1$$

Proof. (Theorem 3.17)

Remark that the l_q -rNSP implies the l_1 -rNSP and the l_p -rNSP ($p \leq q$) in the forms

$$\|\mathbf{v}_S\|_1 \leq \rho\|\mathbf{v}_{\bar{S}}\|_1 + \tau s^{1-1/q}\|A\mathbf{v}\|, \quad (3.5)$$

$$\|\mathbf{v}_S\|_p \leq \frac{\rho}{s^{1-1/q}}\|\mathbf{v}_{\bar{S}}\|_1 + \tau s^{1/p-1/q}\|A\mathbf{v}\|, \quad (3.6)$$

for all $\mathbf{v} \in \mathbb{C}^N$ and all $S \subset [N]$ with $\text{card}(S) \leq s$. In view of inequality (2.5), applying Theorem 2.16 with S chosen as an index set of s largest non-negative components of \mathbf{x} we get

$$\|\mathbf{z} - \mathbf{x}\|_1 \leq \frac{1 + \rho}{1 - \rho}(\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\sigma_s(\mathbf{x})_1) + \frac{2\tau}{1 - \rho}s^{1-1/q}\|A(\mathbf{z} - \mathbf{x})\|. \quad (3.7)$$

Next, we choose S as an index set of s largest non-negative components of $\mathbf{z} - \mathbf{x}$, and apply Lemma 2.18 to it, which yields

$$\|\mathbf{z} - \mathbf{x}\|_p \leq \|(\mathbf{z} - \mathbf{x})_{\bar{S}}\|_p + \|(\mathbf{z} - \mathbf{x})_S\|_p \leq \frac{1}{s^{1-1/p}} \|\mathbf{z} - \mathbf{x}\|_1 + \|(\mathbf{z} - \mathbf{x})_S\|_p.$$

In view of inequality (2.6), we get

$$\begin{aligned} \|\mathbf{z} - \mathbf{x}\|_p &\leq \frac{1}{s^{1-1/p}} \|\mathbf{z} - \mathbf{x}\|_1 + \frac{\rho}{s^{1-1/p}} \|(\mathbf{z} - \mathbf{x})_{\bar{S}}\|_1 + \tau s^{1/p-1/q} \|A(\mathbf{z} - \mathbf{x})\| \\ &\leq \frac{1+\rho}{s^{1-1/p}} \|\mathbf{z} - \mathbf{x}\|_1 + \tau s^{1/p-1/q} \|A(\mathbf{z} - \mathbf{x})\|. \end{aligned} \quad (3.8)$$

Substituting (2.7) into the above inequality (2.8), gives us

$$\begin{aligned} \|\mathbf{z} - \mathbf{x}\|_p &\leq \frac{(1+\rho)^2}{(1-\rho)} \frac{1}{s^{1-1/p}} (\|\mathbf{z}\|_1 - \|\mathbf{x}\|_1 + 2\sigma_s(\mathbf{x})_1) \\ &\quad + \frac{(3+\rho)}{(1-\rho)} \tau s^{1/p-1/q} \|A(\mathbf{z} - \mathbf{x})\|. \end{aligned}$$

Setting $C := (1+\rho)^2/(1-\rho)$ and $D := (3+\rho)\tau/(1-\rho)$ and substituting for it above gives us the desired inequality from Theorem 2.17. ■

3.5 Matrix Analysis: Coherence

It's not easy to verify if a certain matrix is suitable for Basis Pursuit, that is, to check if the given matrix has the NSP. Thus it would be practical if we could find an easy way to calculate a quantity from which we can guarantee the success of the recovery algorithm. The *coherence* is such a quantity.

Definition 3.19. Let $A \in \mathbb{C}^{m \times N}$ be a matrix with l_2 -normalized columns, $\mathbf{a}_1, \dots, \mathbf{a}_N$ such that $\|\mathbf{a}_i\|_2 = 1$ for all $i \in [N]$. The *coherence* $\mu = \mu(A)$ of the matrix A is defined as

$$\mu := \max_{1 \leq i \neq j \leq N} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|$$

The l_1 -coherence μ_1 is defined for $s \in [N-1]$ as

$$\mu_1(s) := \max_{i \in [N]} \max \left\{ \sum_{j \in S} |\langle \mathbf{a}_i, \mathbf{a}_j \rangle|, S \subset [N], \text{card}(S) = s, i \notin S \right\}$$

From the definition of coherence we have the following result:

Theorem 3.20. Let $A \in \mathbb{C}^{m \times N}$ be a matrix with l_2 -normalized columns. If

$$\mu_1(s) + \mu_1(s-1) < 1,$$

then every s -sparse vector $\mathbf{x} \in \mathbb{C}^N$ is exactly recovered from the vector $\mathbf{y} = A\mathbf{x}$ via Basis Pursuit.

3. Compressed Sensing

Proof. If we can show that for a given matrix $A \in \mathbb{C}^{m \times N}$, with l_2 -normalized columns satisfying $\mu_1(s) + \mu_1(s-1) < 1$ also possess the NSP of order s , then according to Theorem 2.10, every s -sparse vector \mathbf{x} recovered via Basis Pursuit, is unique.

Let $\mathbf{a}_1, \dots, \mathbf{a}_N$ denote the columns of A . If $\mathbf{v} \in \ker(A)$, then $A\mathbf{v} = \mathbf{0}$ can be rewritten as the sum of the column vectors of A as the sum $\sum_{j=1}^N v_j \mathbf{a}_j = \mathbf{0}$. Since the $\mu_1(s)$ coherence function is a statement about the sum of the inner products of pairwise columns of A , it will be convenient to rewrite a particular weight v_i as:

$$v_i = v_i \langle \mathbf{a}_i, \mathbf{a}_i \rangle = - \sum_{j=1, j \neq i}^N v_j \langle \mathbf{a}_j, \mathbf{a}_i \rangle$$

Splitting $[N]$ into two parts S and \bar{S} and taking the absolute value of v_i it follows that

$$|v_i| \leq \sum_{l \in \bar{S}} |v_l| |\langle \mathbf{a}_l, \mathbf{a}_i \rangle| + \sum_{j \in S, j \neq i} |v_j| |\langle \mathbf{a}_j, \mathbf{a}_i \rangle|$$

Summing over all $i \in S$ we get

$$\begin{aligned} \|\mathbf{v}_S\|_1 &= \sum_{i \in S} |v_i| \leq \sum_{l \in \bar{S}} |v_l| \sum_{i \in S} |\langle \mathbf{a}_l, \mathbf{a}_i \rangle| + \sum_{j \in S} |v_j| \sum_{i \in S, i \neq j} |\langle \mathbf{a}_j, \mathbf{a}_i \rangle| \\ &\leq \sum_{l \in \bar{S}} |v_l| \mu_1(s) + \sum_{j \in S} |v_j| \mu_1(s-1) = \mu_1(s) \|\mathbf{v}_{\bar{S}}\|_1 + \mu_1(s-1) \|\mathbf{v}_S\|_1 \end{aligned}$$

Rearranging and neglecting the terms between the first term and the last in the equations above we get

$$(1 - \mu_1(s-1)) \|\mathbf{v}_S\|_1 \leq \mu_1(s) \|\mathbf{v}_{\bar{S}}\|_1$$

using the first part of the condition in the theorem, $\mu_1(s) + \mu_1(s-1) < 1$, we get the desired inequality

$$\|\mathbf{v}_S\|_1 < \|\mathbf{v}_{\bar{S}}\|_1$$

which concludes the proof. ■

CHAPTER 4

Neural Networks And Deep Learning

Neural Networks (NNs), a biologically-inspired programming paradigm, provide among the best solutions to many problems where the neural network has a lot of training data. Consequently, NNs perform well in areas such as image classification, speech recognition, natural language processing, and so forth. However, NNs have reportedly been seen to suffer from instabilities in image reconstruction. Our goal is to get mathematical insight to why these instabilities occur. Most of the material presented in this chapter is based on the upcoming book [7]. We begin by introducing the basic elements of NNs, in the framework of supervised machine learning.

4.1 Supervised Machine Learning

In supervised machine learning, the objective is to learn a mapping f from an input space U to an output space V . The 2-tuple set $\{(u_i, v_i)\}_{i=1}^n$, $u_i \in U$, $v_i \in V$ is referred to as the training set. By varying how we define the output space V , we can choose what our function f should model.

For instance, if v_i takes the values -1 or 1, then f may model a binary classification problem or if v_i takes three or more discrete values, then f may model a multi classification problem.

Image reconstruction can be cast as a supervised learning problem as following : Let $U \subseteq \mathbb{R}^m$ and $V \subseteq \mathbb{R}^N$ such that $\{(\mathbf{y}_i, \mathbf{x}_i)\}_{i=1}^n$ is the training set with \mathbf{x}_i as images and \mathbf{y}_i the corresponding measurements. Our goal is then to find a matrix A such that we can express

$$\mathbf{y}_i = A\mathbf{x}_i + \epsilon, \quad i \in \{1, \dots, n\}$$

which corresponds to learning the mapping $f(\mathbf{u}_i) = \mathbf{v}_i + \epsilon$, $i \in \{1, \dots, n\}$. It is desirable that this mapping has a small training error, i.e. $f(\mathbf{u}_i) \approx \mathbf{v}_i$ for $i = 1, \dots, n$, but also that the mapping generalises well to new data $\{(\mathbf{y}_k, \mathbf{x}_k)\}$ which is close to, but not a part of the training set.

The task of computing the mapping f is known as training the model, such a function could for instance be a NN modelling a certain problem.

4. Neural Networks And Deep Learning

4.2 Design

The term *Neural Networks*, covers a large class of models and learning methods. Here we try to give a description that is sufficient for the later analysis.

Definition 4.1. Let $n_0, \dots, n_L \in \mathbb{N}$. Let $W_l : \mathbb{R}^{n_{l-1}} \rightarrow \mathbb{R}^{n_l}$ for $l = 1, 2, \dots, L$ be affine maps. Let $\rho_1, \dots, \rho_{L-1} : \mathbb{R} \rightarrow \mathbb{R}$ be activation functions and let $n_0 = M$ and $n_L = N$. Then a map $\Psi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ given by

$$\Psi(\mathbf{y}) = W_L(\rho_{L-1}(\dots \rho_1(W_1(\mathbf{y}))\dots))$$

is called a *Neural Network*.

The type of network defined above, is called a feedforward network. The name emphasizes that the network takes an input \mathbf{y} and feeds it forward through the network by an alternating sequence of affine maps and non-linear activation functions. It is common to visualize the network as a graph, where the nodes of the graph are the artificial neurons and the entries of W_l are the weights assigned to each edge of the graph. At a given layer, each neuron takes in a sum of linear combination from the outputs of the neurons in the previous layer, applies a bias, applies an activation function and finally feeds it forward to the next layer of neurons.

The *architecture* characterizes the network and consists of the following:

- The depth L , given by the number of layers,
- The width n_0, \dots, n_L of each layer
- The choice of activation function σ .

The number of parameters in a given network grows rather quickly as a function of the number of layers and the width of each layer. More specifically, the number of parameters in a NN is given by:

$$d = (n_1 \times n_0 + n_2 \times n_1 + \dots + n_L \times n_{L-1}) + (n_1 + n_2 + \dots + n_L) \quad (4.1)$$

where n_0, \dots, n_L denotes the number of neurons per layer and $0, \dots, L$ is the number of layers.

The choice of activation function may be important to the performance of the network. Common activation functions in the literature include:

$$\begin{aligned} ReLU(x) &= \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} \\ \tanh(x) &= \frac{e^{2x} - 1}{e^{2x} + 1} \\ \sigma(x) &= \frac{1}{1 + e^{-x}} \end{aligned}$$

If we look closer at the sigmoid function and the hyperbolic tangent, we observe that for large negative input or large positive input, the derivative of these functions will be close to zero. This is known as the *vanishing gradient problem*, which can cause the training of the network to be slow, which is the topic of the next section.

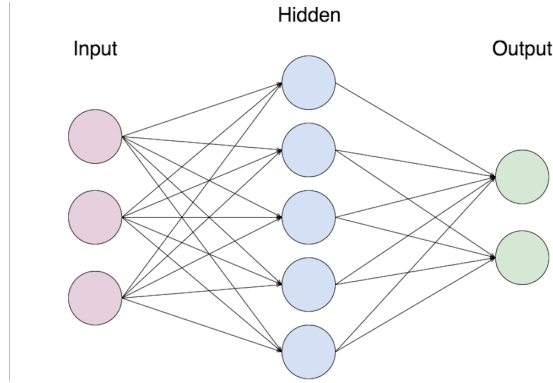


Figure 4.1: Extremely simplified NN with a single hidden layer.

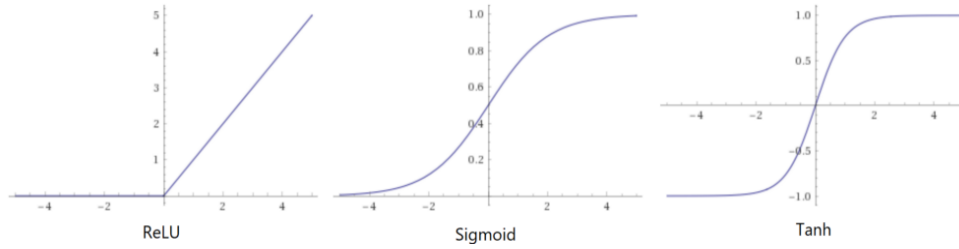


Figure 4.2: Different choices of activations functions.

4.3 Training

The training of a NN is to find values for the weights and biases such that the NN performs well on a given data set. To simplify notation, we will let θ be defined as the set of all trainable parameters. For a neural network this will typically be:

$$\theta = \{\mathbf{W}_1, \dots, \mathbf{W}_L, \mathbf{b}_1, \dots, \mathbf{b}_L\}$$

We write Ψ_θ to emphasize that this NN uses θ as its parameters.

One way to measure how well the network is trained, is by evaluating its performance on a given set by a *cost function*.

Cost function

Several cost functions may be defined on the network. A popular choice is the *quadratic* cost function, also known as the *mean squared error (MSE)* is given by:

$$C(\theta | (\mathbf{u}_1, \mathbf{v}_1), \dots, (\mathbf{u}_n, \mathbf{v}_n)) = \frac{1}{2n} \sum_{i=1}^n \|\Psi_\theta(\mathbf{u}_i) - \mathbf{v}_i\|^2 \quad (4.2)$$

here n is the total number of samples, \mathbf{u}_i the input sample into the network and \mathbf{v}_i the corresponding validation sample.

Next, we want to minimize the cost function, we do this by recasting the

4. Neural Networks And Deep Learning

problem into an optimization problem, where try to find values for all the trainable parameters d (recall 4.1), that is to attempt to solve:

$$\min_{\theta \in \mathbb{R}^d} C(\theta) \quad (4.3)$$

This problem is nonlinear, nonconvex and there is generally no possibility of computing global minimizers of it. Even worse, since a NN needs a huge number of parameters to perform well and a large set of training data, it is a major computational challenge.

Gradient Descent Methods

To make the training of the NN's practical regarding the computational demand, a first-order optimization method is usually used to solve (4.3). These are methods based on the gradient descent. In a gradient descent method we pick a starting point θ_0 , then we follow the path of steepest descent, that is, we produce the sequence:

$$\theta_{i+1} = \theta_i - \tau_i \nabla C(\theta_i), \quad i = 0, 1, \dots \quad (4.4)$$

Calculating the gradient requires one pass through all the training data. When the amount of training data is large, this ceases to be computationally feasible. However, there is at least one way to reduce the number of computational steps, known as *Stochastic Gradient Descent*. There are several variations of the method, they mostly differ by the way they set a step length τ in each iteration, but the main idea they have in common, is to compute the gradient from a smaller subset of the training data and thereby reducing the computation needed at step i . The common algorithm for computing the gradient ∇C , is known as *backpropagation*.

4.4 The Universal Approximation Theorem

In supervised machine learning we are assuming that there exists an underlying "ground truth" function that we are approximating. The Universal Approximation Theorem gives justification for the use of NNs to approximate such a function.

Theorem 4.2. (*Universal Approximation Theorem*) Let $\sigma \in \mathcal{C}(\mathbb{R})$ and let $K \subset \mathbb{R}^n$ be compact.

Then for any $g \in \mathcal{C}(\mathbb{R}^n)$ and any $\epsilon > 0$, there exists a set of parameters $k \in \mathbb{N}$, $c_1, \dots, c_k \in \mathbb{R}$, $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^n$, $b_1, \dots, b_k \in \mathbb{R}$ such that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ defined as

$$f(\mathbf{x}) = \sum_{i=1}^k c_i \sigma(\mathbf{w}_i^T \mathbf{x} + b_i) \quad (4.5)$$

satisfies

$$\max_{\mathbf{x} \in K} |f(\mathbf{x}) - g(\mathbf{x})| < \epsilon$$

for all $\mathbf{x} \in \mathbb{R}^n$ if and only if σ is not a polynomial.

In other words, the class of real-valued neural networks with one hidden layer are dense in $\mathcal{C}(\mathbb{R}^n)$, in the topology of uniform convergence on compact sets, if and only if σ is not a polynomial.

We may easily see that the activation function σ can not possibly be a polynomial if the above theorem is to hold. Consider a polynomial σ of degree m , then for every choice of $\mathbf{x} \in \mathbb{R}^n$ and $b \in \mathbb{R}$, $\sigma(\mathbf{w}^T \mathbf{x} + b)$ is polynomial of total degree at most m , and thus does not span $\mathcal{C}(\mathcal{S})$.

The Universal Approximation Theorem states that the class of NNs with one hidden layer is already extremely rich. However, there are reasons for considering deeper NNs, where the network architecture consist of several layers, since they might approximate certain functions more efficiently than shallow NNs, i.e. networks with few layers.

4.5 The Success Of Deep Learning

Deep Learning, which is a subfield in Artificial Neural Networks, has at least been around since the 1960s [9]. However, it wasn't until DL became practically feasible that it started to get wide-spread attention.

NNs has achieved a great number of successes during the last decade, yielding state-of-the-art performance on a range of challenging machine learning problems. Since the list of successes in different problems is quite long, we will only point out two examples that are milestones in history of Deep Learning applied to image classification tasks.

The first milestone was achieved in 2012, when a Deep Learning-based classifier named *AlexNet* won the *ImageNet Large Scale Visual Recognition Challenge (ILSVRC)*. Prior to 2012, the winners were not Deep Learning based. Not only did it win, but it also achieved an error reduction of nearly 10% over the previous year's winner. Since then, all winners of the ILSVRC contest have used deep learning. In less than a decade, the error rate was reduced from 30% to less than 4%. In 2015, Deep Learning reached a second milestone, when it reduced the image classification error rate to sub 5%, which is the error rate incurred by humans. As a result of this performance, the term 'superhuman' has been used to describe Deep Learning's performance on image classification.

To the authors knowledge, deep learning is now the method of choice for most image classification tasks.

4.6 Instabilities In Deep Learning

Even though Deep Learning has had great success in image classification problems, its performance might also be unstable, particularly against small perturbations in their input. Given a trained classifier C and an input image \mathbf{x}_0 , it is possible to find a visually imperceptible perturbation \mathbf{r} for which the network misclassify the image $\mathbf{x}_0 + \mathbf{r}$ such that

$$C(\mathbf{x}_0) \neq C(\mathbf{x}_0 + \mathbf{r})$$

There even exist methods for finding perturbation that purposefully achieves the misclassification effect. One such method is known as *DeepFool*, see [7]. Around 2014, DeepFool was used to find a small perturbation, that was applied to the previously mentioned *ImageNet* dataset of images, such that the perturbation achieved a classification error rate of over 75 %. Making matters even worse, the perturbation are universal and transferrable, meaning that the same perturbation can be applied to a set of images, and consequently fool, i.e. misclassify, several NNs.

Why do instabilities occur?

With the existence of visually imperceptible changes to an image that fools a NN one may wonder what Deep Learning actually learns. It would be reasonable to assume that humans perceive the abstract characteristics of everyday life images, in a way that is more resilient to imperceptible changes, compared to the current versions of NNs. Some insight might be gain from the so-called *false structure* phenomenon. We will give an example illustrating the phenomenon in shallow NNs:

Consider the set of 3×3 images with a light stripe that is either horizontal or vertical, illustrated in Figure 4.3 taken from [7]. In the vertical case, the light stripe takes the value $1 + a$ and the dark stripes take the value $1 - a$ for some $0 \leq a \leq 1/4$. In the horizontal case, the light stripe takes the value $1 - a$ and the dark stripes take the value $-a$. Next, define the binary labelling function

$$C_0(\mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{x} \text{ has a horizontal light stripe,} \\ -1 & \text{if } \mathbf{x} \text{ has a vertical light stripe.} \end{cases}$$

Now consider a one-layer NN which tries to approximate the above function C_0 defined by

$$C(\mathbf{x}) = \begin{cases} +1 & \text{if the sum of the pixel values are } \leq 3, \\ -1 & \text{if the sum of the pixel values are } > 3. \end{cases}$$

By the way it is designed, this NN has a high success rate and only fails in the case where $a = 0$. However, a small perturbation to the pixel values of the input image can cause C to give the wrong answer, even though the underlying

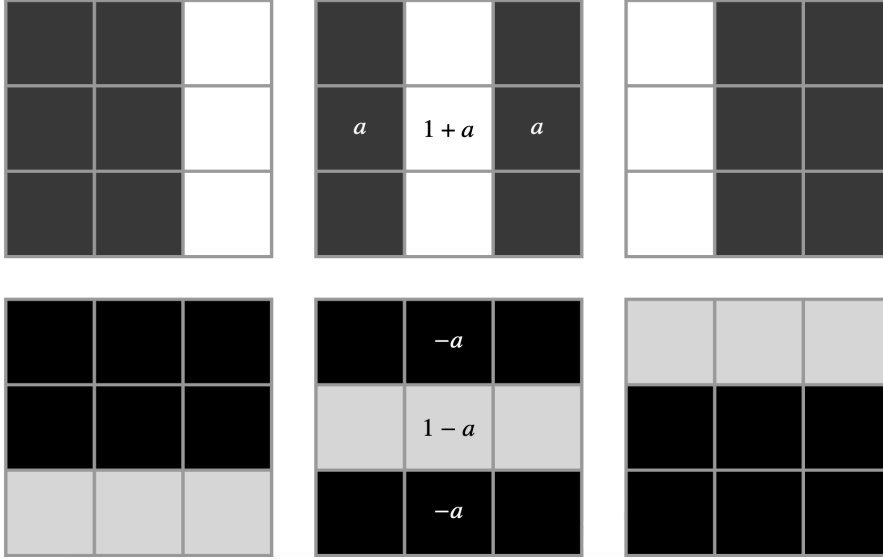


Figure 4.3: 3×3 images with either horizontal or vertical stripes.

ground truth function C_0 is stable regarding such a perturbation. Making the perturbation small enough, it can be invisible to the human eye.

This simple example illustrates that even though the NN may be successful, it does not capture the true underlying structure that defines the function, namely, the stripe. Thus it has learned a false structure and has become unstable in the face of small perturbations.

One might think that the shallow network defined above is too simple in order to capture the underlying structure, as opposed to a deeper network. However, [7] shows that even when one considers families of fully connected, feedforward NNs with ReLu activation function, there are uncountably many classification functions (such as the example above), that will cause trained NNs to be unstable. Unfortunately, the false structure can never be proven with standard mathematical tools, in the sense that as long as NNs are created with inaccurate computations the true minimizers are rarely, if ever found [10].

4.7 Consequences of the false structure phenomenon

In 2019, Thesing, Antun and Hansen published an article where they state the false structure conjecture [10]. In the paper, they provide negative and positive consequences that we shall state here for the sake of completeness.

Negative consequences:

- i) The success of DL in classification is not due to networks learning the structures that humans do, thus instability can never be removed until one guarantees that no false structure is learned.
- ii) Since one does not know which structure the network has learned, it becomes hard to conclude what DL learns and thus difficult to understand its classifications.

4. Neural Networks And Deep Learning

Positive consequences:

- i) Deep learning captures structures that humans do not, thus the structure may have information that the human might not capture. This structure could be useful if characterized properly.

Having knowledge of the false structure phenomenon, might aid in constructing better and more stable NNs in the future.

4.8 Reasons To Go Beyond Compressed Sensing

With the success of Deep Learning in image classification problems and others, one may want to know how well NNs would perform on inverse problems such as image reconstruction compared to more well established techniques such as Compressed Sensing. Other than trying to find the upper limit on the accuracy of NNs in image reconstruction from the fewest possible measurements, there are other reasons for exploring the methods of Deep Learning. There are at least three reasons:

- i) Given a model, such as sparsity and compressibility, Compressed Sensing requires one to handcraft a reconstruction procedure, such as l_1 -minimization, that exploits this structure. It may be challenging to craft an effective reconstruction procedure for an arbitrary image model, therefore a method that is simpler to apply might be advantageous.
- ii) Optimization-based recovery procedures might be slow.
- iii) Finding the optimal parameters of an iterative optimization solver is delicate and can have a significant effect on reconstruction quality.

Once the NN is trained for image reconstruction, it could be faster than iterative optimization-based procedures, since reconstruction only requires a single forward propagation through the network. Also, the time expensive computations may be done at a more convenient time, than optimization-methods.

4.9 Deep Learning for inverse problems

With the aim of understanding the limits of Deep Learning in inverse problems, we end this chapter with a general description of how NNs may solve the discrete, linear inverse problem.

Consider the problem of recovering an unknown vector $\mathbf{x} \in \mathbb{C}^N$ from measurements $\mathbf{y} = A\mathbf{x} + \epsilon \in \mathbb{C}^m$, where $A \in \mathbb{C}^{m \times N}$. Even though most of the examples presented in this thesis have been from signal and image reconstruction, the theoretical results investigated should also apply to the abstract setting. To apply Deep Learning to this problem, we assume there is a training set

$$\mathcal{T} = (\mathbf{y}_i, \mathbf{x}_i)_{i=1}^K \subset \mathbb{C}^m \times \mathbb{C}^N, \mathbf{y}_i = A\mathbf{x}_i,$$

4.9. Deep Learning for inverse problems

consisting of images \mathbf{x}_i and their measurements \mathbf{y}_i . The goal is to feed this data to the NN for it to implicitly learn a reconstruction map $\Psi : \mathbb{C}^m \rightarrow \mathbb{C}^N$ that performs well on the images of interest.

CHAPTER 5

Instabilities In Deep Learning For Inverse Problems

5.1 Universal Instabilities In Inverse Reconstruction Methods

Suppose we have trained a NN to approximate an inverse mapping, $\Psi : \mathbb{C}^m \rightarrow \mathbb{C}^N$ where $m \ll N$. For the NN to be useful, it should perform well on the training set and test set. However, if the NN recovers two vectors \mathbf{x}, \mathbf{x}' whose difference lies in the null space of A , then there exist a ball around $\mathbf{y} = A\mathbf{x}$, where the NN becomes unstable with respect to small perturbations to the input \mathbf{y} . Before we state this consequence as a theorem, recall that the local ϵ -Lipschitz constant of a function ϕ at $\mathbf{y} \in \mathbb{C}^m$, is defined as

$$L^\epsilon(\phi, \mathbf{y}) = \sup_{0 < d_2(\mathbf{z}, \mathbf{y}) \leq \epsilon} \frac{d_1(\phi(\mathbf{z}), \phi(\mathbf{y}))}{d_2(\mathbf{z}, \mathbf{y})}$$

Theorem 5.1. (*Universal Instability Theorem*) Let d_1 and d_2 be metrics on \mathbb{C}^N and \mathbb{C}^m respectively, $A : \mathbb{C}^N \rightarrow \mathbb{C}^m$ a linear map, and $\Psi : \mathbb{C}^m \rightarrow \mathbb{C}^N$ a continuous reconstruction map. Suppose there exist $\mathbf{x}, \mathbf{x}' \in \mathbb{C}^N$ and $\eta > 0$ such that

$$d_1(\Psi(A\mathbf{x}), \mathbf{x}) < \eta, \quad d_1(\Psi(A\mathbf{x}'), \mathbf{x}') < \eta, \quad (5.1)$$

$$d_2(A\mathbf{x}, A\mathbf{x}') \leq \eta. \quad (5.2)$$

Then there is a closed non-empty ball $\mathcal{B} \subset \mathbb{C}^m$ centred at $\mathbf{y} = A\mathbf{x}$ such that the local ϵ -Lipschitz constant at any $\tilde{\mathbf{y}} \in \mathcal{B}$ is bounded from below:

$$L^\epsilon(\Psi, \tilde{\mathbf{y}}) \geq \frac{1}{\epsilon}(d_1(\mathbf{x}, \mathbf{x}') - 2\eta), \quad \epsilon \geq \eta. \quad (5.3)$$

Proof. By the definition of the supremum we have that

$$L^\epsilon(\Psi, \mathbf{y}) = \sup_{0 < d_2(\mathbf{z}, \mathbf{y}) \leq \epsilon} \frac{d_1(\Psi(\mathbf{z}), \Psi(\mathbf{y}))}{d_2(\mathbf{z}, \mathbf{y})} \geq \frac{d_1(\Psi(A\mathbf{x}), \Psi(A\mathbf{x}'))}{d_2(A\mathbf{x}, A\mathbf{x}')}$$

5. Instabilities In Deep Learning For Inverse Problems

Applying the reverse triangle inequality twice we get

$$\frac{d_1(\Psi(A\mathbf{x}), \Psi(A\mathbf{x}'))}{d_2(A\mathbf{x}, A\mathbf{x}')} \geq \frac{d_1(\mathbf{x}, \Psi(A\mathbf{x}')) - d_1(\Psi(A\mathbf{x}), \mathbf{x})}{d_2(A\mathbf{x}, A\mathbf{x}')} \geq \frac{d_1(\mathbf{x}, \mathbf{x}') - d_1(\Psi(A\mathbf{x}), \mathbf{x}) - d_1(\Psi(A\mathbf{x}'), \mathbf{x}')}{d_2(A\mathbf{x}, A\mathbf{x}')}$$

Applying assumption (5.1) and making the bound sharper with $\epsilon \geq \eta$, we get

$$\frac{d_1(\mathbf{x}, \mathbf{x}') - d_1(\Psi(A\mathbf{x}), \mathbf{x}) - d_1(\Psi(A\mathbf{x}'), \mathbf{x}')}{d_2(A\mathbf{x}, A\mathbf{x}')} > \frac{d_1(\mathbf{x}, \mathbf{x}') - 2\eta}{\epsilon}$$

Neglecting the terms in between we get

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

Next, let $\eta_1 = d_1(\Psi(A\mathbf{x}), \mathbf{x})$, and observe that $\eta_1 < \eta$ by (5.1).

Since Ψ is continuous at $\mathbf{y} = A\mathbf{x}$, there exist a $\delta > 0$ such that for all $\tilde{\mathbf{y}} \in \mathbb{C}^m$ with $d_2(\mathbf{y}, \tilde{\mathbf{y}}) < \delta$, we have that $d_1(\Psi(\mathbf{y}), \Psi(\tilde{\mathbf{y}})) \leq \eta - \eta_1$. Thus, for a specific $\delta > 0$, we get that $d_1(\Psi(\mathbf{y}), \Psi(\tilde{\mathbf{y}})) \leq \eta - \eta_1$ for all $\tilde{\mathbf{y}} \in \mathcal{B}_1$, where $\mathcal{B}_1 = \{\tilde{\mathbf{y}} \in \mathbb{C}^m : d_2(\mathbf{y}, \tilde{\mathbf{y}}) < \delta\}$. Next, let $\mathcal{B}_2 = \{\tilde{\mathbf{y}} \in \mathbb{C}^m : d_2(A\mathbf{x}', \tilde{\mathbf{y}}) \leq \eta\}$ and $\mathcal{B} = \mathcal{B}_1 \cap \mathcal{B}_2$. Observe that $\mathcal{B} \neq \emptyset$, since $\mathbf{y} \in \mathcal{B}$. Thus, for any $\tilde{\mathbf{y}} \in \mathcal{B}$ we have

$$\begin{aligned} L^\epsilon(\Psi, \tilde{\mathbf{y}}) &= \sup_{0 < d_2(\mathbf{z}, \tilde{\mathbf{y}}) \leq \epsilon} \frac{d_1(\Psi(\mathbf{z}), \Psi(\tilde{\mathbf{y}}))}{d_2(\mathbf{z}, \tilde{\mathbf{y}})} \geq \frac{d_1(\Psi(A\mathbf{x}'), \Psi(\tilde{\mathbf{y}}))}{d_2(A\mathbf{x}', \tilde{\mathbf{y}})} \\ &\geq \frac{d_1(\Psi(A\mathbf{x}), \Psi(A\mathbf{x}')) - d_1(\Psi(A\mathbf{x}), \Psi(\tilde{\mathbf{y}}))}{d_2(A\mathbf{x}', \tilde{\mathbf{y}})} \\ &\geq \frac{d_1(\mathbf{x}, \mathbf{x}') - d_1(\Psi(A\mathbf{x}), \mathbf{x}) - d_1(\Psi(A\mathbf{x}'), \mathbf{x}') - d_1(\Psi(A\mathbf{x}), \Psi(\tilde{\mathbf{y}}))}{d_2(A\mathbf{x}', A\mathbf{x})} \\ &\geq \frac{d_1(\mathbf{x}, \mathbf{x}') - \eta_1 - \eta - (\eta - \eta_1)}{d_2(A\mathbf{x}', A\mathbf{x})} \\ &\geq \frac{1}{\epsilon}(d_1(\mathbf{x}, \mathbf{x}') - 2\eta). \end{aligned}$$

which is the desired inequality (5.3). ■

5.2 False Positives

Under the same assumptions, as stated in The Universal Instability, we have the following proposition:

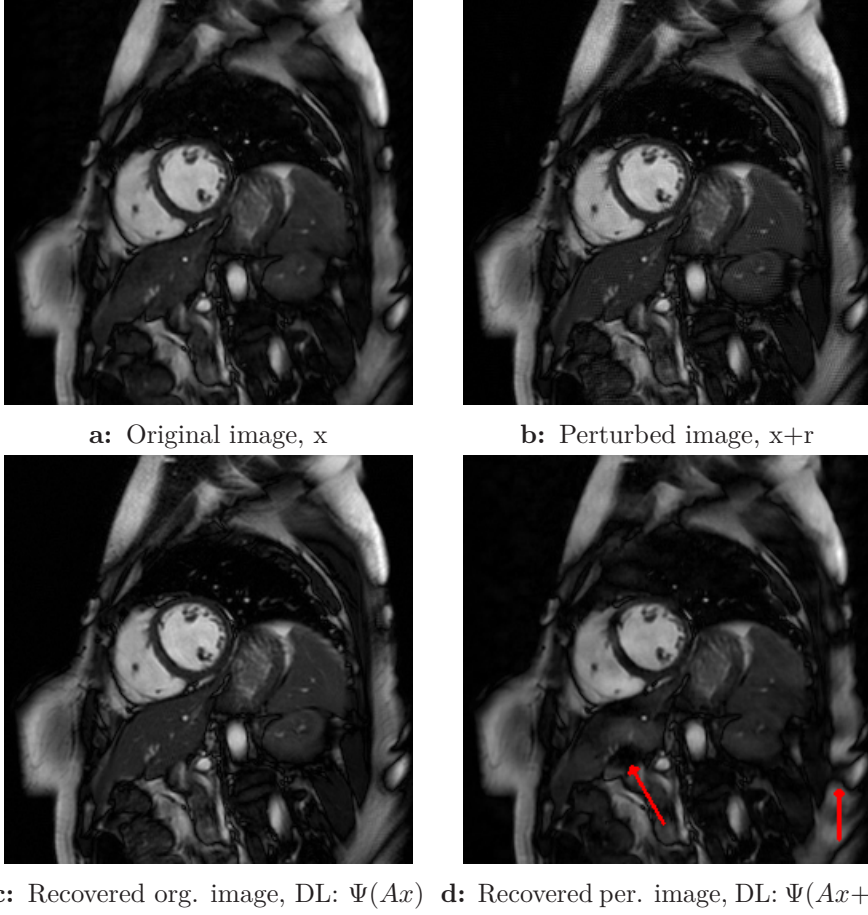


Figure 5.1: False Positive, images from [11].

Proposition 5.2. (*False Positives*) Suppose the conditions from Theorem 5.1 hold, and that the metrics are translation invariant. Then there exists vectors $\mathbf{z} \in \mathbb{C}^N$, $\mathbf{e} \in \mathbb{C}^m$ with $d_1(\mathbf{0}, \mathbf{z}) \geq d_1(\mathbf{x}, \mathbf{x}')$, $d_2(\mathbf{0}, \mathbf{e}) \leq \eta$, and closed non-empty balls $\mathcal{B}_{\mathbf{x}}, \mathcal{B}_{\mathbf{e}}, \mathcal{B}_{\mathbf{z}}$ centred at $\mathbf{x}, \mathbf{e}, \mathbf{z}$ respectively such that

$$d_1(\Psi(A\tilde{\mathbf{x}} + \tilde{\mathbf{e}}), \tilde{\mathbf{x}} + \tilde{\mathbf{z}}) \leq \eta, \quad (5.4)$$

for all $\tilde{\mathbf{x}} \in \mathcal{B}_{\mathbf{x}}, \tilde{\mathbf{e}} \in \mathcal{B}_{\mathbf{e}}, \tilde{\mathbf{z}} \in \mathcal{B}_{\mathbf{z}}$.

Proof. That the metrics d_1, d_2 are translation invariant, means that $d_1(\mathbf{u} + \mathbf{w}, \mathbf{v} + \mathbf{w}) = d_1(\mathbf{u}, \mathbf{v})$ and $d_2(\mathbf{a} + \mathbf{c}, \mathbf{b} + \mathbf{c}) = d_2(\mathbf{a}, \mathbf{b})$ for all $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{C}^N$ and $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{C}^m$ respectively.

Now for proving proposition 5.2, let $\mathbf{z} = \mathbf{x}' - \mathbf{x}$ and $\mathbf{e} = A(\mathbf{x}' - \mathbf{x})$. Translation invariance, implies that $d_1(\mathbf{0}, \mathbf{z}) = d_1(\mathbf{x}', \mathbf{x})$ and by assumption (5.2) in Theorem 5.1, we get $d_2(A\mathbf{x}', A\mathbf{x}) = d_2(\mathbf{0}, \mathbf{e}) \leq \eta$.

5. Instabilities In Deep Learning For Inverse Problems

From assumption 5.1 in Theorem 5.1 and substituting the defined vectors, we get

$$d_1(\Psi(A\mathbf{x} + \mathbf{e}), \mathbf{x} + \mathbf{z}) = d_1(\Psi(A\mathbf{x}'), \mathbf{x}') < \eta.$$

By applying the continuity of Ψ , we arrive at inequality (5.4). \blacksquare

Proposition 5.2 tells us that if the NN recovers two images \mathbf{x}, \mathbf{x}' well (Assumption 5.1), whose difference $\mathbf{x} - \mathbf{x}'$ lies close to the null space of A (Assumption 5.2), then there exist a small perturbation \mathbf{e} , which can be added to the input measurements, such that the network approximately ($< \eta$) recovers an image with an additional element \mathbf{z} , resulting in the image $\mathbf{x} + \mathbf{z}$. This type of instability may have serious consequences, for instance in medical imaging. A *false positive* is a detail, for example a tumour, which is not present in the original image \mathbf{x} , but is present in the reconstructed image. This is illustrated in Figure 5.1. As such, if the image were to be used for diagnosis, then the consequences for the patient might be unfortunate.

The troubling matter of these false positives have also been noted at the 2020 Facebook fastMRI challenge [12] : *"Such hallucinatory features are not acceptable and especially problematic if they mimic normal structures that are either not present or actually abnormal. [...] Neural network models can be unstable as demonstrated via adversarial perturbation studies."*

5.3 False Negatives

Similarly, under the same assumptions as stated in The Universal Instability, we also have the following result:

Proposition 5.3. (*False negatives*) Suppose the conditions from Theorem 5.1 hold, and that the metrics are translation invariant. Then there exists vectors $\mathbf{z} \in \mathbb{C}^N, \mathbf{e} \in \mathbb{C}^m$ with $d_1(\mathbf{0}, \mathbf{z}) \geq d_1(\mathbf{x}, \mathbf{x}'), d_2(\mathbf{0}, \mathbf{e}) \leq \eta$, and closed non-empty balls $\mathcal{B}_{\mathbf{x}}, \mathcal{B}_{\mathbf{e}}, \mathcal{B}_{\mathbf{z}}$ centred at $\mathbf{x}, \mathbf{e}, \mathbf{z}$ respectively such that

$$d_1(\Psi(A(\tilde{\mathbf{x}} + \tilde{\mathbf{z}}) + \tilde{\mathbf{e}}), \tilde{\mathbf{x}}) \leq \eta, \quad (5.5)$$

for all $\tilde{\mathbf{x}} \in \mathcal{B}_{\mathbf{x}}, \tilde{\mathbf{e}} \in \mathcal{B}_{\mathbf{e}}, \tilde{\mathbf{z}} \in \mathcal{B}_{\mathbf{z}}$.

Proof. The proof is almost identical to the proof of Proposition 5.2, except we let $\mathbf{e} = A(\mathbf{x} - \mathbf{x}')$ which gives

$$\begin{aligned} d_1(\Psi(A(\mathbf{x}' + \mathbf{z}) + \mathbf{e}), \mathbf{x}') &= d_1(\Psi(A(\mathbf{x}' + (\mathbf{x}' - \mathbf{x}) + A(\mathbf{x} - \mathbf{x}')), \mathbf{x}') \\ &= d_1(\Psi(A\mathbf{x}'), \mathbf{x}') < \eta. \end{aligned}$$

\blacksquare

Proposition 5.3 tells us that there exist a perturbation \mathbf{e} , such that the reconstruction method Ψ , fails to extract the vector \mathbf{z} from the original image $\mathbf{x} + \mathbf{z}$. Thus a *false negative* is a detail that is present in the original image, but is washed out by the network. For instance in Figure 5.2, the network fails to recover the additional details of the bird and the SIAM logo. [11] States that

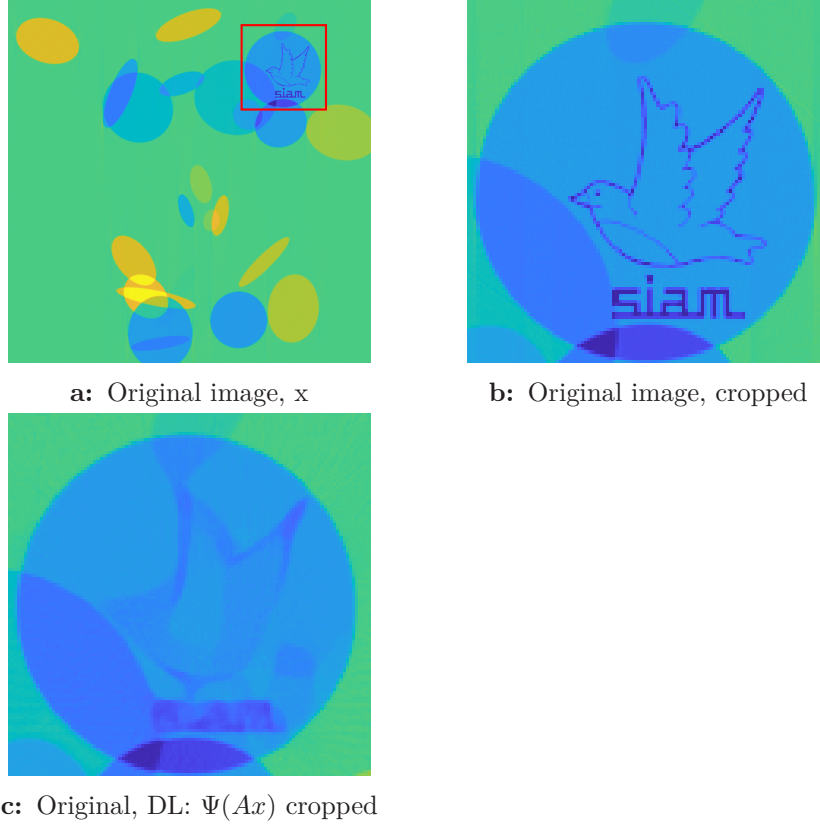


Figure 5.2: False Negatives, images from [11].

the networked was trained to recover ellipses and that the details were not a part of the training set, and thus one may disregard this case. However, it does demonstrate that there may be details of importance that may be washed out.

The existence of false positives and false negatives of reconstruction methods lacking kernel awareness, demonstrates the importance of stability and robustness. This might particularly have an impact on Artificial Intelligence (AI) areas, where Deep Learning techniques are designed to perform tasks that has previously been performed by humans. The European Commission has also taken an interest [13]: *"On AI, trust is a must, not a nice to have. [...] The new AI regulation will make sure that Europeans can trust what AI has to offer. [...] High-risk AI systems will be subject to strict obligations before they can be put on the market: [requiring] High level of robustness, security and accuracy."*

5.4 Stability Through Kernel Awareness

A consequence of the rNSP is that no s -sparse vector can be close to the null space of A . Thus the case, as in assumption 5.2 ($d_2(Ax, Ax') \leq \eta$) in the Universal Instability Theorem, should not occur in CS techniques where x and x' are s -sparse vectors. If the vectors are only close to s -sparse and there is some noise in the measurement vector, Theorem 3.17 tells us that any perturbation in

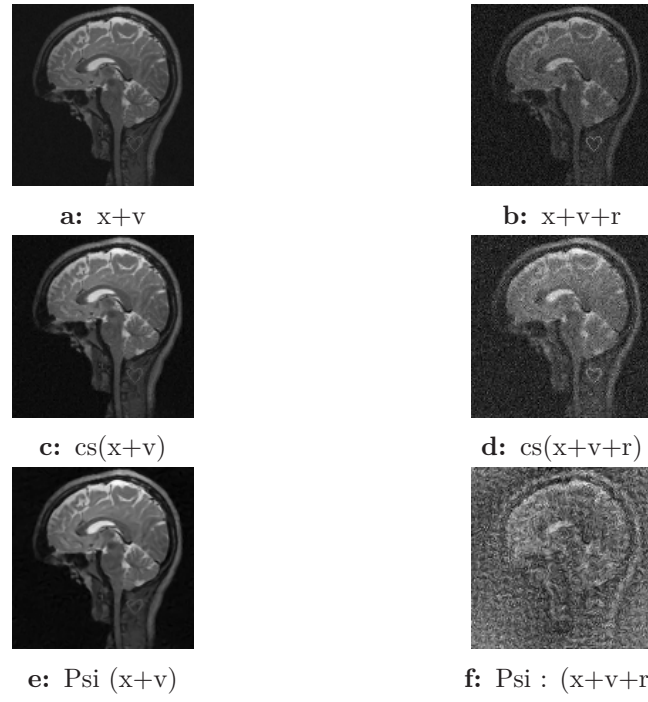


Figure 5.3: CS stable recovery vs unstable DL recovery, figures from [11].

y of magnitude η , yields an error in the recovery of x of magnitude a constant times η and how accurate the recovery is. The stability and accuracy of CS techniques is illustrated in Figure 5.4.



a: Original image, cropped



b: Original, DL: $\Psi(Ax)$ cropped



c: Original image, CS cropped x

Figure 5.4: CS stable recovery vs unstable DL recovery, figures from [11].

Appendices

Bibliography

- [1] Higham, N. J., *Accuracy and Stability of Numerical Algorithms*, 2nd ed. SIAM, 2002.
- [2] Lohne, M., “Parseval reconstruction networks,” M.S. thesis, University of Oslo, 2019.
- [3] Holand, I., “The loss of the sleipner condeep platform,” in *DIANA Computational Mechanics ‘94*, Springer Netherlands, 1994.
- [4] Skeel, R., “Roundoff error and the patriot missile,” *SIAM News*, vol. Volume 25, no. Number 4, page 11, Jul. 1992.
- [5] Inquiry Report, B. of, “Ariane 5, flight 501 failure,” European Space Agency, Tech. Rep., 1996.
- [6] Foucart, S. and Rauhut, H., *A Mathematical Introduction to Compressive Sensing*. Springer, 2013.
- [7] Adcock, B. and Hansen, A., *Structured Compressed Sensing, Imaging and Learning*. Cambridge University Press, coming soon.
- [8] Yoshida, T. and Ohki, K., “Natural images are reliably represented by sparse and variable populations of neurons in visual cortex,” *Nat Commun*, vol. 11, 2020.
- [9] Schmidhuber, J., “Deep learning in neural networks: An overview,” vol. 61, Jan. 2015.
- [10] Thesing, L., Antun, V., and Hansen, A. C., “What do ai algorithms actually learn? - on false structures in deep learning,” Jun. 2019.
- [11] Gottschling, N., Antun, V., Adcock, B., and Hansen, A., “The troublesome kernel: Why deep learning for inverse problems is typically unstable,” Jan. 2020.
- [12] Muckley, M. J., Riemenschneider, B., Radmanesh, A., Kim, S., Jeong, G., Ko, J., Jun, Y., Shin, H., Hwang, D., Mostapha, M., Arberet, S., Nickel, D., Ramzi, Z., Ciuciu, P., Starck, J.-L., Teuwen, J., Karkalousos, D., Zhang, C., Sriram, A., Huang, Z., Yakubova, N., Lui, Y., and Knoll, F., “Results of the 2020 fastmri challenge for machine learning mr image reconstruction,” Dec. 2020.

Bibliography

- [13] Commission, E., *Europe fit for the digital age: [Https://digital-strategy.ec.europa.eu/en/news/europe-fit-digital-age-commission-proposes-new-rules-and-actions-excellence-and-trust-artificial](https://digital-strategy.ec.europa.eu/en/news/europe-fit-digital-age-commission-proposes-new-rules-and-actions-excellence-and-trust-artificial)*, Press release, Apr. 2021.