

Inverse Problems in Imaging

Lecture 2

Simon R. Arridge¹

¹Department of Computer Science, University College London, UK

UCL, Term 2

1 Introduction

2 Basic Concepts

- Some examples
- Inverse Problems
- Singular Value Decomposition (SVD)
- Singular Value Decomposition for Operators
- The Continuous-Discrete Case

1 Introduction

2 Basic Concepts

- Some examples
- Inverse Problems
- Singular Value Decomposition (SVD)
- Singular Value Decomposition for Operators
- The Continuous-Discrete Case

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.
- From a computational point of view, we are trying to fit a model to data, which is done by *optimisation*.

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.
- From a computational point of view, we are trying to fit a model to data, which is done by *optimisation*.
- Not all inverse problems are solved by optimisation (but many are).

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.
- From a computational point of view, we are trying to fit a model to data, which is done by *optimisation*.
- Not all inverse problems are solved by optimisation (but many are).
- Not all optimisation problems are about inverse problems. We will do some simple “functional” examples to introduce useful algorithms.

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.
- From a computational point of view, we are trying to fit a model to data, which is done by *optimisation*.
- Not all inverse problems are solved by optimisation (but many are).
- Not all optimisation problems are about inverse problems. We will do some simple “functional” examples to introduce useful algorithms.
- Not all inverse problems involve images. We will do some “toy” problems to get the concepts across.

Inverse Problems in Imaging

Introduction

- Inverse problems involves the inference of the parameters of a model from data. In this course we are considering (mainly) *Image Reconstruction* - i.e. the inferred parameters form an image, but the data comes from many different types of experimental observations.
- From a computational point of view, we are trying to fit a model to data, which is done by *optimisation*.
- Not all inverse problems are solved by optimisation (but many are).
- Not all optimisation problems are about inverse problems. We will do some simple “functional” examples to introduce useful algorithms.
- Not all inverse problems involve images. We will do some “toy” problems to get the concepts across.
- Mostly, we concentrate on *Inverse Problems in Imaging* : i.e. model inversion whose solutions are images.

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$.

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$. We might also have time variation (sometimes called 4D).

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$. We might also have time variation (sometimes called 4D). Hyperspectral images add another dimension. 5D imaging !

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$. We might also have time variation (sometimes called 4D). Hyperspectral images add another dimension. 5D imaging !
- Images are *spatially organised* : the connectivity between pixels is due to physical proximity. We can make use of *structural coherence*, e.g. edges. In time-varying imaging we can exploit *temporal coherence*. In Hyperspectral imaging we can exploit *spectral coherence*.

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$. We might also have time variation (sometimes called 4D). Hyperspectral images add another dimension. 5D imaging !
- Images are *spatially organised* : the connectivity between pixels is due to physical proximity. We can make use of *structural coherence*, e.g. edges. In time-varying imaging we can exploit *temporal coherence*. In Hyperspectral imaging we can exploit *spectral coherence*.
- In other words, imaging comes with a lot of detailed *prior information* of a specific kind. This helps a lot !

Inverse Problems in Imaging

Introduction

What's special about images rather than generic inverse problems ?

- Images are very large : several million pixels in camera images.
- In tomography in the real world we deal with 3D data. Images might be $> 10^8$. We might also have time variation (sometimes called 4D). Hyperspectral images add another dimension. 5D imaging !
- Images are *spatially organised* : the connectivity between pixels is due to physical proximity. We can make use of *structural coherence*, e.g. edges. In time-varying imaging we can exploit *temporal coherence*. In Hyperspectral imaging we can exploit *spectral coherence*.
- In other words, imaging comes with a lot of detailed *prior information* of a specific kind. This helps a lot !
- Compare this to machine learning where the connectivity of a graph representation of data has to be defined in more abstract ways, e.g. using kernels.

Inverse Problems in Imaging

Introduction : Optimisation

Broadly speaking Optimisation requires us to consider four things

Inverse Problems in Imaging

Introduction : Optimisation

Broadly speaking Optimisation requires us to consider four things

- A *metric*. This addresses the question “Optimisation *of* what?”. A metric is a single real positive number that acts as a “figure of merit” of our success at optimisation.

Inverse Problems in Imaging

Introduction : Optimisation

Broadly speaking Optimisation requires us to consider four things

- A *metric*. This addresses the question “Optimisation *of* what?”. A metric is a single real positive number that acts as a “figure of merit” of our success at optimisation.
- A *model*. This addresses the question “what is the metric applied to?”. A model is a mathematical or computational process that predicts a result which maybe a function, vector, list or some other quantity.

Inverse Problems in Imaging

Introduction : Optimisation

Broadly speaking Optimisation requires us to consider four things

- A *metric*. This addresses the question “Optimisation *of* what?”. A metric is a single real positive number that acts as a “figure of merit” of our success at optimisation.
- A *model*. This addresses the question “what is the metric applied to?”. A model is a mathematical or computational process that predicts a result which maybe a function, vector, list or some other quantity.
- *Parameters*. This addresses the question “Optimisation *over* what ?”. These are the variable controls of the model. The aim of optimisation is to adjust the parameters until the figure of metric is in some way “as good as we can get”

Inverse Problems in Imaging

Introduction : Optimisation

Broadly speaking Optimisation requires us to consider four things

- A *metric*. This addresses the question “Optimisation *of* what?”. A metric is a single real positive number that acts as a “figure of merit” of our success at optimisation.
- A *model*. This addresses the question “what is the metric applied to?”. A model is a mathematical or computational process that predicts a result which maybe a function, vector, list or some other quantity.
- *Parameters*. This addresses the question “Optimisation *over* what ?”. These are the variable controls of the model. The aim of optimisation is to adjust the parameters until the figure of metric is in some way “as good as we can get”
- An *algorithm*. This addresses the question “How do we optimise?”

Inverse Problems in Imaging

Introduction : Optimisation

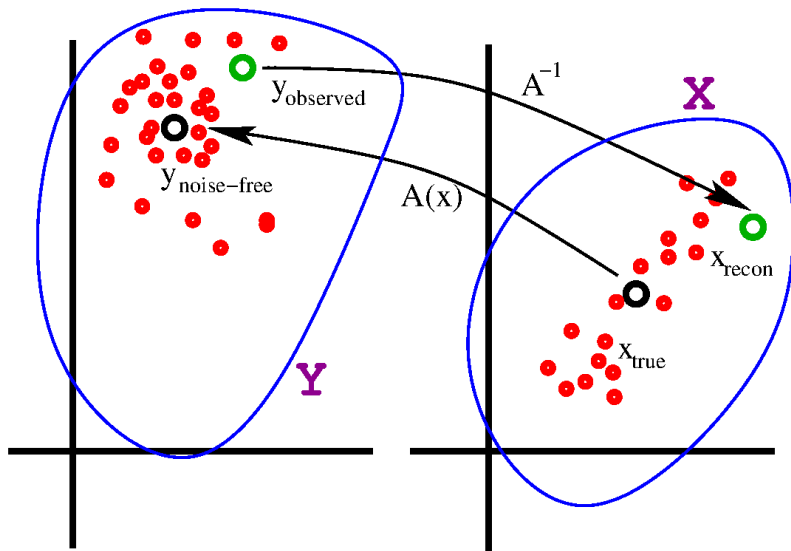
Broadly speaking Optimisation requires us to consider four things

- A *metric*. This addresses the question “Optimisation *of* what?”. A metric is a single real positive number that acts as a “figure of merit” of our success at optimisation.
- A *model*. This addresses the question “what is the metric applied to?”. A model is a mathematical or computational process that predicts a result which maybe a function, vector, list or some other quantity.
- *Parameters*. This addresses the question “Optimisation *over* what ?”. These are the variable controls of the model. The aim of optimisation is to adjust the parameters until the figure of metric is in some way “as good as we can get”
- An *algorithm*. This addresses the question “How do we optimise?”

And for inverse problems we need a fifth thing : Data !

Inverse Problems in Imaging

Introduction : overall picture



Inverse Problems in Imaging

Introduction

Key concepts :

- A space Y , (“Data Space”), where the data lives. Here the “distance” between points is based on the statistics of the measurement noise.
- A space X , (“Solution Space”), where the solution (i.e. the images) live. Here we have different concepts of the “distance” between images.

Commonly used

- 1 Peak Signal to Noise Ratio (PSNR)
- 2 Structural Similarity Index Measure (SSIM)

There are other choices.

- X and Y might be the same. E.g. Image Deblurring, Image Denoising, Inpainting etc. These are *Image Processing* problems.
- In other cases Y is quite different, e.g. in tomography, geophysical imaging etc.

Inverse Problems in Imaging

Introduction : the Forward Problem

At the heart of solving inverse problems is to understand the *Forward Problem*

$$y = A(x)$$

A may be linear or non-linear.

A takes an element of solution space X (i.e. an image) and predicts the corresponding element of data space Y .

- In the case of **denoising** A is just the identity operator
- In the case of **deblurring** A is a blurring operator (stationary or non-stationary)
- In the case of **inpainting** A is a subsampling operator (like a mask)
- In the case of **X-ray tomography** A represents integration along lines.

A major part of the field of inverse problems is to formulate the forward problem correctly, and then to develop a computational model for it. For example in Geophysics imaging it involves solving a Partial Differential Equation in 3D.

The forward problem may often be the computational bottleneck in solving inverse problems.

Inverse Problems in Imaging

Introduction : the Inverse Problem

We want to go the other way : find x from the data y

$$x_{\text{recon}} = A^{-1}(y_{\text{observed}})$$

That's the hard part !

A simple idea : discretise x and y into vectors, then A should be a matrix. Then invert this. We call this the “Brute Force” method!

- But matrix inversion could be very costly. In fact even holding a matrix representation of the operator in memory could be intractable. E.g. 10^6 pixels and 10^6 data makes A of size $10^6 \times 10^6$. About a terabyte.
- We need smarter ideas \Rightarrow “Matrix-Free” methods.
- But sometimes we might do the Brute Force method on a scaled down problem as a kind of “Sanity Check”.

Inverse Problems in Imaging

Introduction : the effect of noise

If we had the ground truth x_{true} and apply the forward model we get *noise-free* data

$$y_{noise-free} = A(x_{true})$$

If we have “solved” the inverse problem, i.e. we have a method implementing A^{-1} , then we will get x_{true} back, obviously.

However this can be misleading. In fact it is sometimes called an *Inverse Crime*. We have neglected the effect of noise. In practice the observations are corrupted by noise (assume this is additive)

$$y_{observed} = A(x_{true}) + \eta$$

The noise η is assumed drawn from a distribution $\eta \sim P(\eta)$.

The reconstruction error $e = x_{recon} - x_{true}$ is then also from some distribution. But what form is it ?

Inverse Problems in Imaging

Introduction : the effect of noise

If we repeat solving the inverse problem with different instantiations of the measurement noise (e.g. just doing different physical experiments, or adding different samples of noise from the model $P(\eta)$), then we will get *samples* of the reconstructed image drawn from an (unknown) *posterior probability distribution* $P_{\text{post}}(x)$.

$$y_{\text{observed}}^{(k)} = A(x_{\text{true}}) + \eta^{(k)} \Leftrightarrow x_{\text{recon}}^{(k)} = A^{-1}(y_{\text{observed}}^{(k)})$$

We might expect that if the noise was well behaved, e.g. Gaussian white noise, then so would be the posterior $P_{\text{post}}(x)$.

However, it usually turns out that posterior distribution of the error is far from being white noise.

The good news is that the mean of $P_{\text{post}}(x)$ is actually the correct solution x_{true} . The bad news is that our confidence estimates of this solution are usually terrible!

The reason is because the forward operator **does not propagate information uniformly**.

Inverse Problems in Imaging

Introduction : Bayesian perspective

From the Bayesian perspective, if we take an observation and invert it, we get the *Maximum Likelihood* solution : i.e. x_{recon} is “most likely” to match the data y_{observed} .

If we take a set of data and reconstructions $\{y_{\text{observed}}^{(k)}, x_{\text{recon}}^{(k)}\}$ then the *Conditional Mean*

$$x_{\text{CM}} = \frac{1}{N} \sum_{k=1}^N x_{\text{recon}}^{(k)}$$

converges to x_{true} in the limit as $N \rightarrow \infty$.

However our estimate of covariance

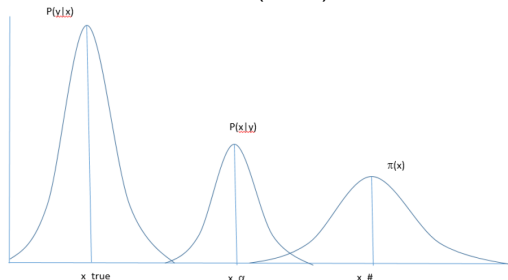
$$C = \frac{1}{N} \sum_{k=1}^N \left(x_{\text{recon}}^{(k)} - x_{\text{CM}} \right) \left(x_{\text{recon}}^{(k)} - x_{\text{CM}} \right)^T$$

Can have eigenvalues tending to ∞

Inverse Problems in Imaging

Introduction : effect of priors

Including instead a *prior distribution* of the solutions $\pi(x)$ produces instead the *Maximum A-Priori* (MAP) estimate.



However the MAP estimate is not equal to x_{true} . It introduces *bias*.

$$x_{\text{bias}} := x_{\text{MAP}} - x_{\text{true}} \neq 0$$

Part of the effort in solving inverse problems is to design (or learn) appropriate priors and determine a tradeoff between bias and variance.

Inverse Problems in Imaging

Introduction : Summary

- Introduced Imaging, Inverse Problems, and Optimisation.
- Inverse Problems involve Forward Problems, which can be the computational bottleneck.
- Inverting a Forward Problem produces the Maximum Likelihood solution. It fits the data but may be unstable with respect to noise.
- Adding prior information can stabilise an inverse problem but introduces bias.

Next we need to analyse these concepts more carefully, and make quantitative statements.

1 Introduction

2 Basic Concepts

- Some examples
- Inverse Problems
- Singular Value Decomposition (SVD)
- Singular Value Decomposition for Operators
- The Continuous-Discrete Case

Basic Concepts

Introduction

- In this section we introduce some key concepts :
 - Non-Uniqueness
 - Non-Existence
 - Instability of solutions
- We'll look at “toy” examples in two dimensions. This will allow us to analyse in a lot of detail and also visualise the ideas as 2D graphics.
- Some of the ideas are “obvious”, but keep in mind how these examples will generalise to large number of dimensions.

Basic Concepts

Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

$$x_1 + 2x_2 = 5 \quad (2.1)$$

Basic Concepts

Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

$$x_1 + 2x_2 = 5 \quad (2.1)$$

Obviously, there is no unique solution because the problem has one equation in two unknowns . The set of possible solutions forms a line on the x_1, x_2 plane.

Basic Concepts

Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

$$x_1 + 2x_2 = 5 \quad (2.1)$$

Obviously, there is no unique solution because the problem has one equation in two unknowns . The set of possible solutions forms a line on the x_1, x_2 plane. Eq.(2.1) corresponds to the matrix equation

$$y = A\mathbf{x}$$

where $A = \begin{pmatrix} 1 & 2 \end{pmatrix}$, $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, $y = 5$. A non-square matrix doesn't have an exact inverse.

Basic Concepts

Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

$$x_1 + 2x_2 = 5 \quad (2.1)$$

Obviously, there is no unique solution because the problem has one equation in two unknowns . The set of possible solutions forms a line on the x_1, x_2 plane. Eq.(2.1) corresponds to the matrix equation

$$y = A\mathbf{x}$$

where $A = \begin{pmatrix} 1 & 2 \end{pmatrix}$, $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, $y = 5$. A non-square matrix doesn't have an exact inverse.

For other values of y , the solutions always lie on lines that are parallel to this case.

Basic Concepts

Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

$$x_1 + 2x_2 = 5 \quad (2.1)$$

Obviously, there is no unique solution because the problem has one equation in two unknowns . The set of possible solutions forms a line on the x_1, x_2 plane. Eq.(2.1) corresponds to the matrix equation

$$y = A \mathbf{x}$$

where $A = \begin{pmatrix} 1 & 2 \end{pmatrix}$, $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, $y = 5$. A non-square matrix doesn't have an exact inverse.

For other values of y , the solutions always lie on lines that are parallel to this case.

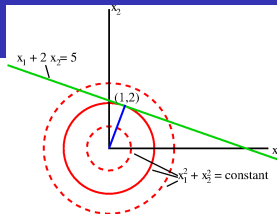
Even if A is square it is not always invertible. So how can we get a useful solution ?

Basic Concepts

Constrained Optimisation

Consider instead the following problem :

$$\begin{array}{ll} \text{minimise} & \Phi = x_1^2 + x_2^2 \\ \text{subject to} & x_1 + 2x_2 = 5 \quad (2.2) \end{array}$$



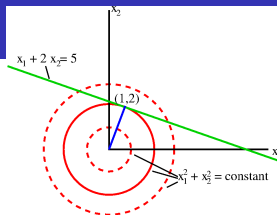
Basic Concepts

Constrained Optimisation

Consider instead the following problem :

$$\begin{array}{ll} \text{minimise} & \Phi = x_1^2 + x_2^2 \\ \text{subject to} & x_1 + 2x_2 = 5 \quad (2.2) \end{array}$$

- The level sets of Φ are circles with value equal to the square of their radius. We can find the solution by starting at a large value $\Phi \sim \infty$ and “shrinking” the circles towards the origin. Obviously the minimum of Φ itself is zero, but when including the constraint we get the solution $(x_1, x_2) = (1, 2)$.

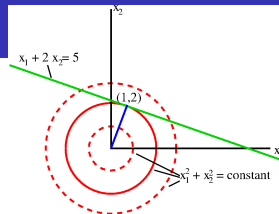


Basic Concepts

Constrained Optimisation

Consider instead the following problem :

$$\begin{array}{ll}\text{minimise} & \Phi = x_1^2 + x_2^2 \\ \text{subject to} & x_1 + 2x_2 = 5 \quad (2.2)\end{array}$$



- The level sets of Φ are circles with value equal to the square of their radius. We can find the solution by starting at a large value $\Phi \sim \infty$ and “shrinking” the circles towards the origin. Obviously the minimum of Φ itself is zero, but when including the constraint we get the solution $(x_1, x_2) = (1, 2)$.
- This is an example of *constrained optimisation*. The unique solution is the one where the surface $\Phi = \text{constant}$ is **tangent to the constraint** $x_1 + 2x_2 = 5$. The expression $\Phi = x_1^2 + x_2^2$ is a *norm* (i.e. a metric, or distance measure) on the Euclidean space \mathbb{R}^2 with basis vectors $\hat{\mathbf{x}}_1 = (1, 0)$, $\hat{\mathbf{x}}_2 = (0, 1)$.

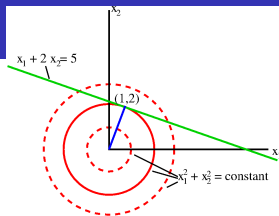
Basic Concepts

Constrained Optimisation

Consider instead the following problem :

$$\begin{aligned} \text{minimise} \quad & \Phi = x_1^2 + x_2^2 \\ \text{subject to} \quad & x_1 + 2x_2 = 5 \quad (2.2) \end{aligned}$$

- The level sets of Φ are circles with value equal to the square of their radius. We can find the solution by starting at a large value $\Phi \sim \infty$ and “shrinking” the circles towards the origin. Obviously the minimum of Φ itself is zero, but when including the constraint we get the solution $(x_1, x_2) = (1, 2)$.
- This is an example of *constrained optimisation*. The unique solution is the one where the surface $\Phi = \text{constant}$ is **tangent to the constraint** $x_1 + 2x_2 = 5$. The expression $\Phi = x_1^2 + x_2^2$ is a *norm* (i.e. a metric, or distance measure) on the Euclidean space \mathbb{R}^2 with basis vectors $\hat{\mathbf{x}}_1 = (1, 0)$, $\hat{\mathbf{x}}_2 = (0, 1)$.
- The solution that we get is the *Minimum Norm* solution. It gives a unique point out of the otherwise infinite possible set of solutions.



Basic Concepts

Moore-Penrose Inverse

Comment : We could also have got the answer using the *Moore-Penrose Inverse* (MPI) of the matrix A .

Basic Concepts

Moore-Penrose Inverse

Comment : We could also have got the answer using the *Moore-Penrose Inverse* (MPI) of the matrix A . The MPI is a way of getting an inverse even if a matrix is non-square. It is defined

$$A^{\dagger} = \begin{cases} A^T(AA^T)^{-1} & \text{if } A \text{ has more columns than rows} \\ (A^TA)^{-1}A^T & \text{if } A \text{ has more rows than columns} \end{cases}.$$

Basic Concepts

Moore-Penrose Inverse

Comment : We could also have got the answer using the *Moore-Penrose Inverse* (MPI) of the matrix A . The MPI is a way of getting an inverse even if a matrix is non-square. It is defined

$$A^{\dagger} = \begin{cases} A^T(AA^T)^{-1} & \text{if } A \text{ has more columns than rows} \\ (A^TA)^{-1}A^T & \text{if } A \text{ has more rows than columns} \end{cases}.$$

In this case we take the first expression and get

$$A^T = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, (AA^T) = (5), (AA^T)^{-1} = \frac{1}{5} \Rightarrow A^{\dagger} = \begin{pmatrix} 1/5 \\ 2/5 \end{pmatrix} \quad (2.3)$$

Applying A^{\dagger} to the “data” $y = 5$ gives the solution $\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$.

Basic Concepts

Moore-Penrose Inverse

Comment : We could also have got the answer using the *Moore-Penrose Inverse* (MPI) of the matrix A . The MPI is a way of getting an inverse even if a matrix is non-square. It is defined

$$A^{\dagger} = \begin{cases} A^T(AA^T)^{-1} & \text{if } A \text{ has more columns than rows} \\ (A^TA)^{-1}A^T & \text{if } A \text{ has more rows than columns} \end{cases}.$$

In this case we take the first expression and get

$$A^T = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, (AA^T) = (5), (AA^T)^{-1} = \frac{1}{5} \Rightarrow A^{\dagger} = \begin{pmatrix} 1/5 \\ 2/5 \end{pmatrix} \quad (2.3)$$

Applying A^{\dagger} to the “data” $y = 5$ gives the solution $\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$.

Note that for any data y the solution always lies along the same line, i.e.

$$\mathbf{x} = \lambda \begin{pmatrix} 1 \\ 2 \end{pmatrix} \text{ with } \lambda = \frac{y}{5}.$$

Basic Concepts

Moore-Penrose Inverse

Comment : We could also have got the answer using the *Moore-Penrose Inverse* (MPI) of the matrix A . The MPI is a way of getting an inverse even if a matrix is non-square. It is defined

$$A^{\dagger} = \begin{cases} A^T(AA^T)^{-1} & \text{if } A \text{ has more columns than rows} \\ (A^T A)^{-1} A^T & \text{if } A \text{ has more rows than columns} \end{cases}.$$

In this case we take the first expression and get

$$A^T = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, (AA^T) = (5), (AA^T)^{-1} = \frac{1}{5} \Rightarrow A^{\dagger} = \begin{pmatrix} 1/5 \\ 2/5 \end{pmatrix} \quad (2.3)$$

Applying A^{\dagger} to the “data” $y = 5$ gives the solution $\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$.

Note that for any data y the solution always lies along the same line, i.e.

$$\mathbf{x} = \lambda \begin{pmatrix} 1 \\ 2 \end{pmatrix} \text{ with } \lambda = \frac{y}{5}.$$

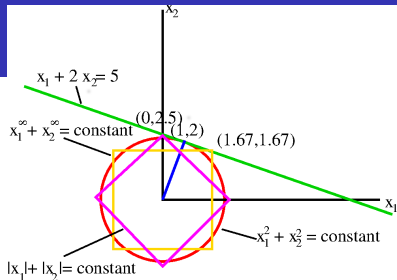
This idea is somewhat obvious. What’s interesting is that we get the solution by minimising the particular *Euclidean norm* with the constraint that the value got by applying the forward model to the solution **agrees with the data**.

Basic Concepts

Generalising the norm

Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \end{array} \quad (2.4)$$



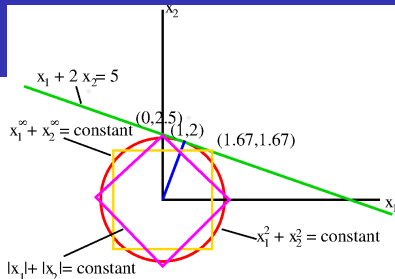
Basic Concepts

Generalising the norm

Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \quad (2.4) \end{array}$$

When $1 < p < \infty$ the surfaces $\Phi = \text{constant}$ are *convex* and solutions can be found relatively easily (at least “by hand”).



Basic Concepts

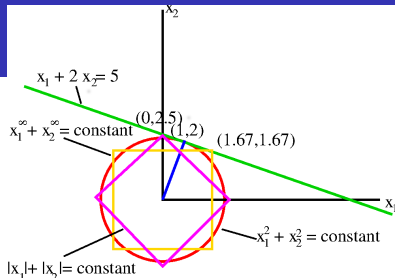
Generalising the norm

Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \end{array} \quad (2.4)$$

When $1 < p < \infty$ the surfaces $\Phi = \text{constant}$ are *convex* and solutions can be found relatively easily (at least “by hand”).

We see that we can obtain any solution between $(0, 2.5)$ at $p = 1$ and $(1\frac{2}{3}, 1\frac{2}{3})$ for $p = \infty$.



Basic Concepts

Generalising the norm

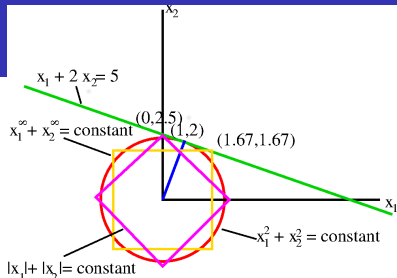
Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \quad (2.4) \end{array}$$

When $1 < p < \infty$ the surfaces $\Phi = \text{constant}$ are *convex* and solutions can be found relatively easily (at least “by hand”).

We see that we can obtain any solution between $(0, 2.5)$ at $p = 1$ and $(1\frac{2}{3}, 1\frac{2}{3})$ for $p = \infty$.

However, actually calculating the solution is not as simple as solving a linear problem like eq. 2.3. The Moore-Penrose Inverse doesn't work in this case.



Basic Concepts

Generalising the norm

Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \end{array} \quad (2.4)$$

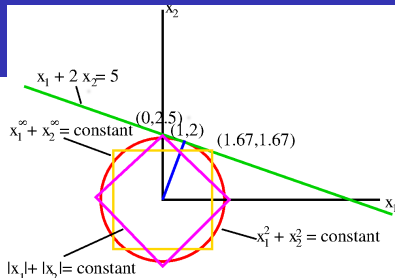
When $1 < p < \infty$ the surfaces $\Phi = \text{constant}$ are *convex* and solutions can be found relatively easily (at least “by hand”).

We see that we can obtain any solution between $(0, 2.5)$ at $p = 1$ and $(1\frac{2}{3}, 1\frac{2}{3})$ for $p = \infty$.

However, actually calculating the solution is not as simple as solving a linear problem like eq. 2.3. The Moore-Penrose Inverse doesn't work in this case.

Notice that for the case $p = 1$ the solution for any equation $ax_1 + bx_2 = c$ will (almost always) lie on the coordinate axes : $\mathbf{x} = \begin{pmatrix} 0 \\ c/b \end{pmatrix}$ or $\mathbf{x} = \begin{pmatrix} c/a \\ 0 \end{pmatrix}$.

I.e. the solution will (almost always) have only one non-zero component.



Basic Concepts

Generalising the norm

Now consider that we take a different norm and ask for the solution to

$$\begin{array}{ll} \text{minimise} & \Phi = |x_1|^p + |x_2|^p \\ \text{subject to} & x_1 + 2x_2 = 5. \end{array} \quad (2.4)$$

When $1 < p < \infty$ the surfaces $\Phi = \text{constant}$ are *convex* and solutions can be found relatively easily (at least “by hand”).

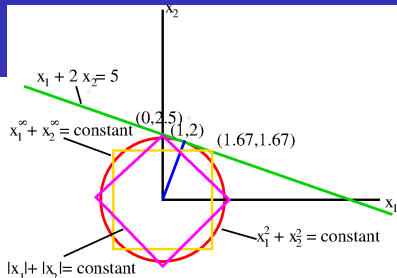
We see that we can obtain any solution between $(0, 2.5)$ at $p = 1$ and $(1\frac{2}{3}, 1\frac{2}{3})$ for $p = \infty$.

However, actually calculating the solution is not as simple as solving a linear problem like eq. 2.3. The Moore-Penrose Inverse doesn't work in this case.

Notice that for the case $p = 1$ the solution for any equation $ax_1 + bx_2 = c$ will (almost always) lie on the coordinate axes : $\mathbf{x} = \begin{pmatrix} 0 \\ c/b \end{pmatrix}$ or $\mathbf{x} = \begin{pmatrix} c/a \\ 0 \end{pmatrix}$.

I.e. the solution will (almost always) have only one non-zero component.

By contrast, as $p \rightarrow \infty$ the solution components tend to be equal $|x_1| \simeq |x_2|$.



Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions. Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case!

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions.

Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case!

The value of any number to the zeroth power is defined

$$x^0 = \begin{cases} 0 & \text{if } |x| = 0 \\ 1 & \text{otherwise} \end{cases}$$

So for $p = 0$ the functional Φ just counts the number of non zeros in \mathbf{x} .

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions.

Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case!

The value of any number to the zeroth power is defined

$$x^0 = \begin{cases} 0 & \text{if } |x| = 0 \\ 1 & \text{otherwise} \end{cases}$$

So for $p = 0$ the functional Φ just counts the number of non zeros in \mathbf{x} .

This means minimising Φ minimises the number of non-zeros in the solution, subject to being consistent with the data.

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions.

Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case!

The value of any number to the zeroth power is defined

$$x^0 = \begin{cases} 0 & \text{if } |x| = 0 \\ 1 & \text{otherwise} \end{cases}$$

So for $p = 0$ the functional Φ just counts the number of non zeros in \mathbf{x} .

This means minimising Φ minimises the number of non-zeros in the solution, subject to being consistent with the data. The solution is called *sparse* when it has few non-zeros

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions. Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case! The value of any number to the zeroth power is defined

$$x^0 = \begin{cases} 0 & \text{if } |x| = 0 \\ 1 & \text{otherwise} \end{cases}$$

So for $p = 0$ the functional Φ just counts the number of non zeros in \mathbf{x} . This means minimising Φ minimises the number of non-zeros in the solution, subject to being consistent with the data. The solution is called *sparse* when it has few non-zeros

But as we see, at least in simple cases, we get the same result using $p = 1$ (or any $p \in [0, 1]$).

Basic Concepts

About Sparsity

Comment about *Sparsity*:

For values $0 < p < 1$ the level sets of Φ are *concave*.

Optimisation of non-convex functions is much harder than convex functions. Nevertheless, we can see that starting with a large value of Φ and shrinking to the origin, the same solutions are obtained as for solving the $p = 1$ case! The value of any number to the zeroth power is defined

$$x^0 = \begin{cases} 0 & \text{if } |x| = 0 \\ 1 & \text{otherwise} \end{cases}$$

So for $p = 0$ the functional Φ just counts the number of non zeros in \mathbf{x} . This means minimising Φ minimises the number of non-zeros in the solution, subject to being consistent with the data. The solution is called *sparse* when it has few non-zeros

But as we see, at least in simple cases, we get the same result using $p = 1$ (or any $p \in [0, 1]$). This motivates the use of the $p = 1$ norm for priors when we want to enforce sparsity.

Basic Concepts

Model Fitting Example

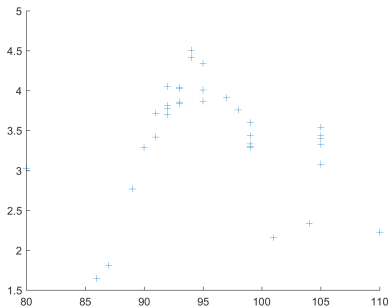
Model Fitting Example

Generate data, $\{x^{(k)}, y^{(k)}; k = 1 \dots N\}$
with the model

$$A_{\text{true}}(x) = 1 + \frac{1}{50}x - \frac{1}{3} \sin\left(\frac{x^2}{30}\right) + e^{\left(-\frac{0.2x^3}{40}\right)}$$

Add 20% multiplicative noise

$$y = (1 + 0.2r)A_{\text{true}}(x), \quad r \sim \mathcal{N}(0, 1).$$



Fit data with model $A_{\text{model}}(x; \xi) = \sum_{m=0}^M \xi_m x^m$ in order to minimise the *Least Squares Error*

$$\Phi = \frac{1}{2} \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|^2 \quad (2.5)$$

Basic Concepts

Model Fitting Example Results

The model $A_{\text{model}}(x; \xi)$ is non-linear in x but linear in ξ .

E.g. for a cubic fit ($M = 4$), Φ is minimised by solving

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix}$$

Basic Concepts

Model Fitting Example Results

The model $A_{\text{model}}(x; \xi)$ is non-linear in x but linear in ξ .

E.g. for a cubic fit ($M = 4$), Φ is minimised by solving

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix}$$

Basic Concepts

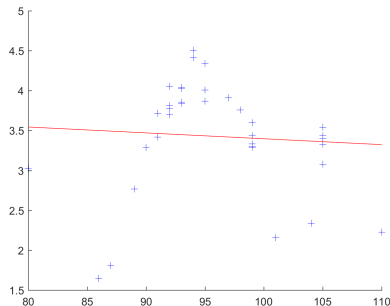
Model Fitting Example Results

The model $A_{\text{model}}(x; \xi)$ is non-linear in x but linear in ξ .

E.g. for a cubic fit ($M = 4$), Φ is minimised by solving

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix}$$

Code example `modelfit.m`
Results $M = 2$ (linear)



Basic Concepts

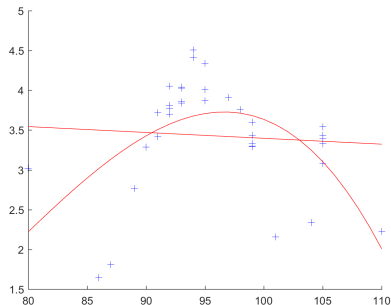
Model Fitting Example Results

The model $A_{\text{model}}(x; \xi)$ is non-linear in x but linear in ξ .

E.g. for a cubic fit ($M = 4$), Φ is minimised by solving

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix}$$

Code example `modelfit.m`
Results $M = 4$ (cubic)



Basic Concepts

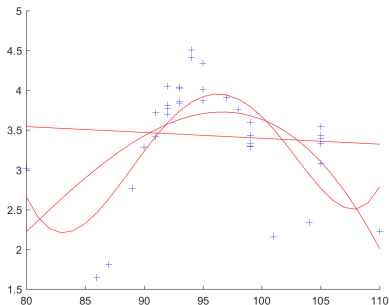
Model Fitting Example Results

The model $A_{\text{model}}(x; \xi)$ is non-linear in x but linear in ξ .

E.g. for a cubic fit ($M = 4$), Φ is minimised by solving

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix}$$

Code example `modelfit.m`
Results $M = 6$ (5th order)

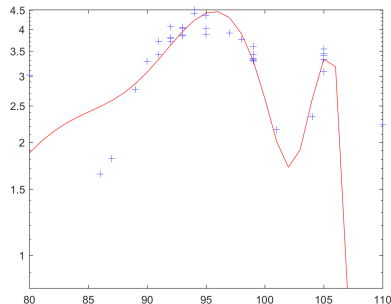


Basic Concepts

Model Fitting Example Results

Increasing the number of fitting parameters leads to *Over-fitting*

Code example `modelfit.m`
Results $M = 36$



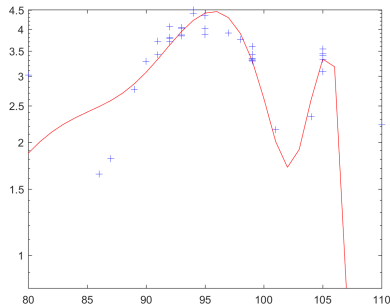
- Closer fitting to the data leads to undesirable features in the solution (noise propagation).

Basic Concepts

Model Fitting Example Results

Increasing the number of fitting parameters leads to *Over-fitting*

Code example `modelfit.m`
Results $M = 36$



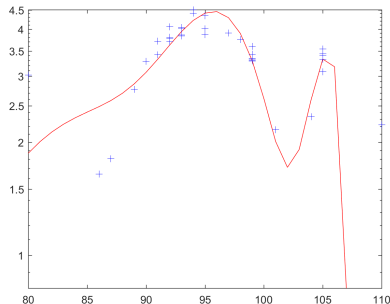
- Closer fitting to the data leads to undesirable features in the solution (noise propagation).
- We should choose the complexity of the model appropriately (number of parameters).

Basic Concepts

Model Fitting Example Results

Increasing the number of fitting parameters leads to *Over-fitting*

Code example `modelfit.m`
Results $M = 36$



- Closer fitting to the data leads to undesirable features in the solution (noise propagation).
- We should choose the complexity of the model appropriately (number of parameters).
- Alternatively : use *Regularisation*. See later !

Basic Concepts

Model Fitting : Comments

- Least-square fitting corresponds to a belief that the noise is i.i.d. Gaussian white noise.

Basic Concepts

Model Fitting : Comments

- Least-square fitting corresponds to a belief that the noise is i.i.d. Gaussian white noise.
- More generally we could define a *data-fitting* functional

$$\Phi = \mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) . \quad (2.6)$$

Basic Concepts

Model Fitting : Comments

- Least-square fitting corresponds to a belief that the noise is i.i.d. Gaussian white noise.
- More generally we could define a *data-fitting* functional

$$\Phi = \mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) . \quad (2.6)$$

Noise model

data fitting functional

Gaussian white

$$\Phi = \frac{1}{2} \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|^2$$

Poissonian

$$\Phi = \sum_{i=k}^N y^{(k)} \log(A_{\text{model}}(x^{(k)}, \xi)) - A_{\text{model}}(x^{(k)}, \xi)$$

Salt and Pepper

$$\Phi = \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|$$

Basic Concepts

Model Fitting : Comments

- Least-square fitting corresponds to a belief that the noise is i.i.d. Gaussian white noise.
- More generally we could define a *data-fitting* functional

$$\Phi = \mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) . \quad (2.6)$$

Noise model

data fitting functional

Gaussian white $\Phi = \frac{1}{2} \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|^2$

Poissonian $\Phi = \sum_{i=k}^N y^{(k)} \log(A_{\text{model}}(x^{(k)}, \xi)) - A_{\text{model}}(x^{(k)}, \xi)$

Salt and Pepper $\Phi = \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|$

- In general, we should analyse (or *learn*) the statistical model of the data and take the negative log of the *likelihood*

$$\mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) = -\log P(\mathbf{y}|A(\mathbf{x}, \xi))$$

Basic Concepts

Model Fitting : Comments

- Least-square fitting corresponds to a belief that the noise is i.i.d. Gaussian white noise.
- More generally we could define a *data-fitting* functional

$$\Phi = \mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) . \quad (2.6)$$

Noise model

data fitting functional

Gaussian white

$$\Phi = \frac{1}{2} \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|^2$$

Poissonian

$$\Phi = \sum_{i=k}^N y^{(k)} \log(A_{\text{model}}(x^{(k)}, \xi)) - A_{\text{model}}(x^{(k)}, \xi)$$

Salt and Pepper

$$\Phi = \sum_{i=k}^N |y^{(k)} - A_{\text{model}}(x^{(k)}, \xi)|$$

- In general, we should analyse (or *learn*) the statistical model of the data and take the negative log of the *likelihood*

$$\mathcal{D}(\mathbf{y}, A(\mathbf{x}, \xi)) = -\log P(\mathbf{y}|A(\mathbf{x}, \xi))$$

- There is no reason to assume that the number of model parameters of the model is the same as the number of data. E.g when reconstructing 3D tomographic images, the number of pixels (voxels) is in some sense “arbitrary”. How many should we choose ?

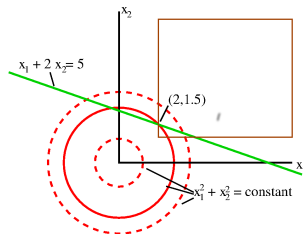
Basic Concepts

Model Fitting : Inequality Constraints

Consider the same problem as before but introduce bound on the solution

$$\begin{array}{ll}\text{minimise} & \Phi = x_1^2 + x_2^2 \\ \text{subject to} & x_1 + 2x_2 = 5 \\ & 2 \leq x_1 \leq 4 \\ & 1 \leq x_2 \leq 5\end{array}$$

which has solution $(x_1, x_2) = (2, 1.5)$,



Unlike eq. 2.2 this cannot be solved by a simple linear method like eq. 2.3. Instead we need *Iterative methods*.

Now we define explicitly the three aspects we mentioned before

Problem 1 : Non-Uniqueness

Solve

$$x_1 + x_2 = 2 \quad (2.7)$$

This problem has an infinite number of solutions. The set of equations are *underdetermined*. To find a single solution we need to add constraints and use a constrained optimisation method.

Basic Concepts

Inverse Problems

Now we define explicitly the three aspects we mentioned before

Problem 1 : Non-Uniqueness

Solve

$$x_1 + x_2 = 2 \quad (2.7)$$

This problem has an infinite number of solutions. The set of equations are *underdetermined*. To find a single solution we need to add constraints and use a constrained optimisation method.

Problem 2 : Non-Existence

Solve

$$\begin{aligned} x_1 &= 1 \\ x_1 &= 3 \end{aligned} \quad (2.8)$$

This problem has no solution. The set of equations are *overdetermined*. We can find a possible solution using unconstrained optimisation.

Basic Concepts

Inverse Problems : Non Existence

Non-Existence is the mirror-image of non-uniqueness.

Basic Concepts

Inverse Problems : Non Existence

Non-Existence is the mirror-image of non-uniqueness.

Problem 2 in matrix form is

$$\begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} x_1$$

Let's use the Moore-Penrose inverse in its overdetermined form :

Basic Concepts

Inverse Problems : Non Existence

Non-Existence is the mirror-image of non-uniqueness.

Problem 2 in matrix form is

$$\begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} x_1$$

Let's use the Moore-Penrose inverse in its overdetermined form : In this case we get

$$A^T = \begin{pmatrix} 1 & 1 \end{pmatrix}, (A^T A) = (2), (A^T A)^{-1} = \frac{1}{2} \Rightarrow A^\dagger = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad (2.9)$$

Basic Concepts

Inverse Problems : Non Existence

Non-Existence is the mirror-image of non-uniqueness.

Problem 2 in matrix form is

$$\begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} x_1$$

Let's use the Moore-Penrose inverse in its overdetermined form : In this case we get

$$A^T = \begin{pmatrix} 1 & 1 \end{pmatrix}, (A^T A) = (2), (A^T A)^{-1} = \frac{1}{2} \Rightarrow A^\dagger = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad (2.9)$$

Applying A^\dagger to the “data” $y = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$ gives the solution $x_1 = 2$.

Basic Concepts

Inverse Problems : Non Existence

Non-Existence is the mirror-image of non-uniqueness.

Problem 2 in matrix form is

$$\begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} x_1$$

Let's use the Moore-Penrose inverse in its overdetermined form : In this case we get

$$A^T = \begin{pmatrix} 1 & 1 \end{pmatrix}, (A^T A) = (2), (A^T A)^{-1} = \frac{1}{2} \Rightarrow A^\dagger = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \end{pmatrix} \quad (2.9)$$

Applying A^\dagger to the “data” $y = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$ gives the solution $x_1 = 2$.

This is the “direct” solution. But it can also be solved by optimisation.

Basic Concepts

Inverse Problems : Non Existence

Consider the following problem :

$$\Phi = (x_1 - 1)^2 + (x_1 - 3)^2 \rightarrow \min \quad (2.10)$$

Basic Concepts

Inverse Problems : Non Existence

Consider the following problem :

$$\Phi = (x_1 - 1)^2 + (x_1 - 3)^2 \rightarrow \min \quad (2.10)$$

It is easy to see we cannot get a solution with $\Phi = 0$ because we can solve a quadratic equation exactly

$$\begin{aligned}(x_1 - 1)^2 + (x_1 - 3)^2 &= 0 \\ \Rightarrow 2x_1^2 - 8x_1 + 10 &= 0 \\ \Rightarrow x_1 &= \frac{4 \pm \sqrt{16 - 20}}{2} \\ &= 2 \pm 2i\end{aligned}$$

i.e. the solution has no real roots.

Basic Concepts

Inverse Problems : Non Existence

Consider the following problem :

$$\Phi = (x_1 - 1)^2 + (x_1 - 3)^2 \rightarrow \min \quad (2.10)$$

It is easy to see we cannot get a solution with $\Phi = 0$ because we can solve a quadratic equation exactly

$$\begin{aligned}(x_1 - 1)^2 + (x_1 - 3)^2 &= 0 \\ \Rightarrow 2x_1^2 - 8x_1 + 10 &= 0 \\ \Rightarrow x_1 &= \frac{4 \pm \sqrt{16 - 20}}{2} \\ &= 2 \pm 2i\end{aligned}$$

i.e. the solution has no real roots.

If we want a real valued solution we can find the minimum :

$$x_1 = 2 \Leftrightarrow \Phi(x_1) = 2$$

This is the *least squares solution*.

Basic Concepts

Non-Uniqueness/Non Existence

Comments

- Non-uniqueness **always occurs** in real life problems because the number of unknowns is “infinite” even if the measurements are finite. We make a problem finite by discretisation, but this implicitly means that we have chosen a particular *basis* for the solution space X .

Comments

- Non-uniqueness **always occurs** in real life problems because the number of unknowns is “infinite” even if the measurements are finite. We make a problem finite by discretisation, but this implicitly means that we have chosen a particular *basis* for the solution space X .
- Non-Existence means that the model cannot match the data exactly. There is some component of the data outside the range of the model.

Comments

- Non-uniqueness **always occurs** in real life problems because the number of unknowns is “infinite” even if the measurements are finite. We make a problem finite by discretisation, but this implicitly means that we have chosen a particular *basis* for the solution space X .
- Non-Existence means that the model cannot match the data exactly. There is some component of the data outside the range of the model.
- Physics tells us that a solution must exist, but our model may be incomplete. This is not measurement noise but *modelling error* (also known as *systematic error*).

Basic Concepts

Non-Uniqueness/Non Existence

Comments

- Non-uniqueness **always occurs** in real life problems because the number of unknowns is “infinite” even if the measurements are finite. We make a problem finite by discretisation, but this implicitly means that we have chosen a particular *basis* for the solution space X .
- Non-Existence means that the model cannot match the data exactly. There is some component of the data outside the range of the model.
- Physics tells us that a solution must exist, but our model may be incomplete. This is not measurement noise but *modelling error* (also known as *systematic error*).
- In practice we can “fix” non-existence by using least-square. I.e. we find the solution that matches the data “as close as possible”. But this may still give problems when modelling error is present.

Basic Concepts

Inverse Problems : Instability of solution to errors

Problem 3 : Discontinuous dependence on data

Consider this problem :

$$A\mathbf{x} = \begin{pmatrix} 1 & 1 + \epsilon \\ 1 - \epsilon & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (2.11)$$

Basic Concepts

Inverse Problems : Instability of solution to errors

Problem 3 : Discontinuous dependence on data

Consider this problem :

$$A\mathbf{x} = \begin{pmatrix} 1 & 1 + \epsilon \\ 1 - \epsilon & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (2.11)$$

The matrix A has determinant ϵ^2 so it has a proper inverse

$$A^{-1} = \frac{1}{\epsilon^2} \begin{pmatrix} 1 & -(1 + \epsilon) \\ -(1 - \epsilon) & 1 \end{pmatrix} \quad (2.12)$$

Thus the inverse grows as an inverse square in the parameter ϵ . This problem has a solution that tends towards instability.

Basic Concepts

Inverse Problems : Instability of solution to errors

Problem 3 : Discontinuous dependence on data

Consider this problem :

$$\mathbf{A}\mathbf{x} = \begin{pmatrix} 1 & 1 + \epsilon \\ 1 - \epsilon & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (2.11)$$

The matrix \mathbf{A} has determinant ϵ^2 so it has a proper inverse

$$\mathbf{A}^{-1} = \frac{1}{\epsilon^2} \begin{pmatrix} 1 & -(1 + \epsilon) \\ -(1 - \epsilon) & 1 \end{pmatrix} \quad (2.12)$$

Thus the inverse grows as an inverse square in the parameter ϵ . This problem has a solution that tends towards instability.

We can do an experiment to see this

Algorithm 1 Statistical trial of ill-posed matrix inversion

for all trials $i = 1 \dots N$ **do**

 draw a random vector of length 2 with Gaussian probability $\eta_i \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right)$

 generate data $b_i = (10, 20) + \eta_i$

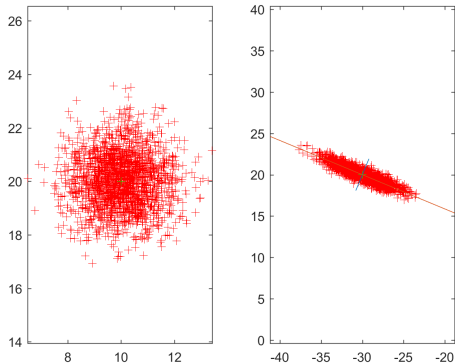
 solve $\mathbf{A}\mathbf{x}_i = b_i$

end for

Basic Concepts

Instability Results

Code example `iptoy.m`



Results for $\epsilon = 1$

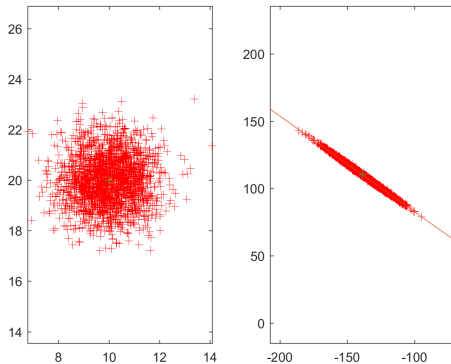
Sample Covariance eigenvalues : $[0.1696, 5.7801]$,

eigenvectors : $\begin{pmatrix} 0.3840 \\ 0.9233 \end{pmatrix}, \begin{pmatrix} -0.9233 \\ 0.3840 \end{pmatrix}$

Basic Concepts

Instability Results

Code example `iptoy.m`



Results for $\epsilon = 0.35$

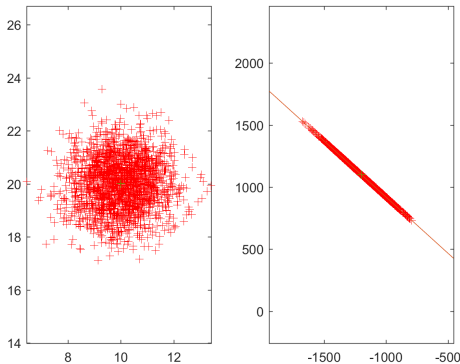
Sample Covariance eigenvalues : $[0.2296, 261.0188]$,

eigenvectors : $\begin{pmatrix} 0.5785 \\ 0.8157 \end{pmatrix}, \begin{pmatrix} -0.8157 \\ 0.5785 \end{pmatrix}$

Basic Concepts

Instability Results

Code example `iptoy.m`



Results for $\epsilon = 0.1$

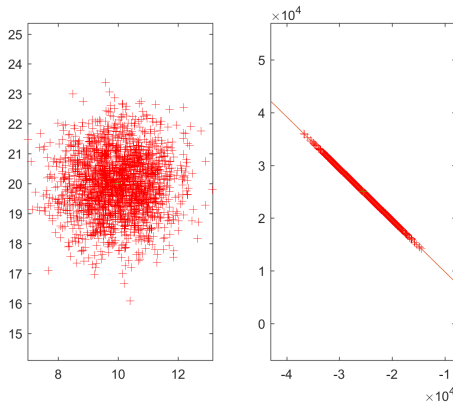
Sample Covariance eigenvalues : $[0.2507, , 38084]$,

eigenvectors : $\begin{pmatrix} 0.6710 \\ 0.7415 \end{pmatrix}, \begin{pmatrix} -0.7415 \\ 0.6710 \end{pmatrix}$

Basic Concepts

Instability Results

Code example `iptoy.m`



Results for $\epsilon = 0.0202$

Sample Covariance eigenvalues : $[0.2556, 2.3806 \times 10^7]$,

eigenvectors : $\begin{pmatrix} 0.6999 \\ 0.7142 \end{pmatrix}, \begin{pmatrix} -0.7142 \\ 0.6999 \end{pmatrix}$

Comments

- As $\epsilon \rightarrow 0$ the distribution of solutions becomes more and more elongated and also rotates clockwise.

Comments

- As $\epsilon \rightarrow 0$ the distribution of solutions becomes more and more elongated and also rotates clockwise.
- Can we determine this behaviour from the matrix A itself ? If so what is the limiting behaviour ?

Comments

- As $\epsilon \rightarrow 0$ the distribution of solutions becomes more and more elongated and also rotates clockwise.
- Can we determine this behaviour from the matrix A itself ? If so what is the limiting behaviour ?
- It turns out we can predict this behaviour accurately using the *Singular Value Decomposition*.

Basic Concepts

Instability Results : Comparison to Learning

Since we have a lot of samples of the data, can we use **learning** to improve the results ?

Basic Concepts

Instability Results : Comparison to Learning

Since we have a lot of samples of the data, can we use **learning** to improve the results ? Results using a simple network and single true data with added noise

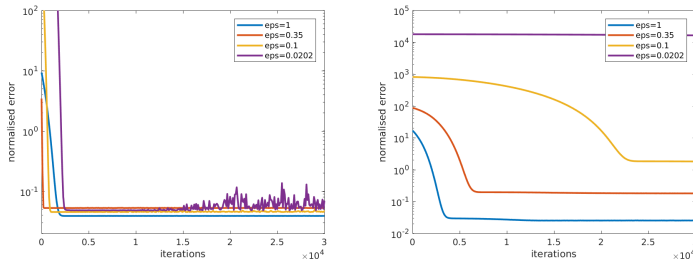


Figure: Convergence of training procedure for forward and inverse network for a test set of 128 samples. Left: forward problem, right: inverse problem

Basic Concepts

Instability Results : Comparison to Learning

Since we have a lot of samples of the data, can we use **learning** to improve the results ? Results using a simple network and single true data with added noise

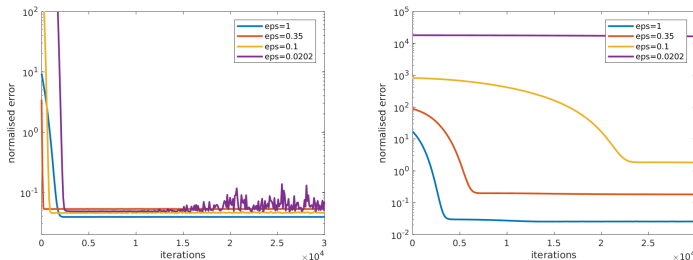


Figure: Convergence of training procedure for forward and inverse network for a test set of 128 samples. Left: forward problem, right: inverse problem

Results show there is a marked difference between the problems!

Basic Concepts

Instability Results : Comparison to Learning

As a more general idea, we also draw the true data from a distribution

Basic Concepts

Instability Results : Comparison to Learning

As a more general idea, we also draw the true data from a distribution

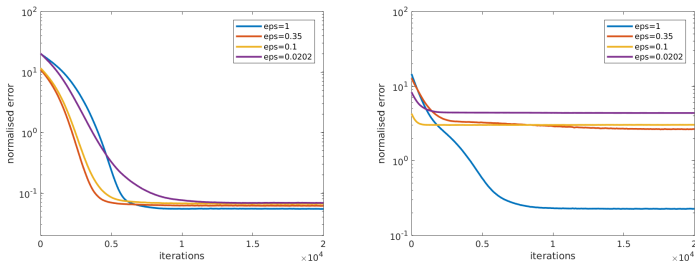


Figure: Convergence of training procedure for forward and inverse network for a test set of 128 samples. Left: forward problem, right: inverse problem

Basic Concepts

Instability Results : Comparison to Learning

As a more general idea, we also draw the true data from a distribution

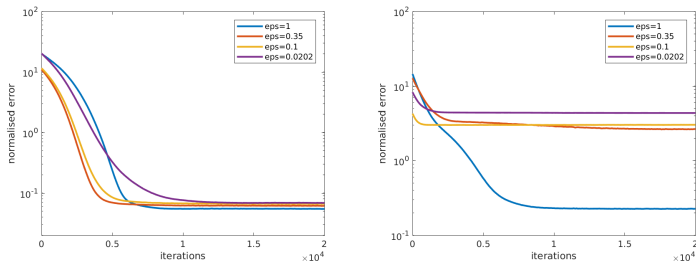


Figure: Convergence of training procedure for forward and inverse network for a test set of 128 samples. Left: forward problem, right: inverse problem

Forward problem is slower to converge, but the inverse problem still is unstable when ϵ is too small

Basic Concepts

Instability Results : Comparison to Learning

As a more general idea, we also draw the true data from a distribution

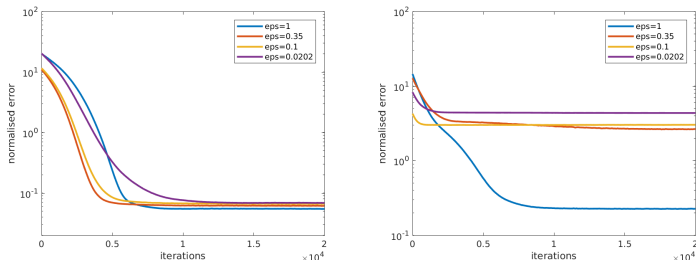


Figure: Convergence of training procedure for forward and inverse network for a test set of 128 samples. Left: forward problem, right: inverse problem

Forward problem is slower to converge, but the inverse problem still is unstable when ϵ is too small

Could different network architectures produce better results ?

Basic Concepts

EigenDecomposition

The eigenvalues of A can be found by solving

$$\lambda^2 - 2\lambda + \epsilon^2 = 0 \quad \Rightarrow \quad \lambda = 1 \pm \sqrt{1 - \epsilon^2}.$$

with corresponding eigenvectors (unnormalised)

$$\mathbf{e}_1 = \begin{pmatrix} \sqrt{1 - \epsilon^2} \\ 1 - \epsilon \end{pmatrix}, \quad \mathbf{e}_2 = \begin{pmatrix} -\sqrt{1 - \epsilon^2} \\ 1 - \epsilon \end{pmatrix}$$

Note: Since A is not symmetric, the eigenvectors are not orthogonal
Taking the limit of these solutions as $\epsilon \rightarrow 0$ we find the limiting values for the eigenvalues :

$$\lambda_1 \rightarrow 2 - \frac{1}{2}\epsilon, \quad \lambda_2 \rightarrow \frac{1}{2}\epsilon$$

with corresponding eigenvectors

$$\mathbf{e}_1 \rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{e}_2 \rightarrow \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

Note: In the limit A is symmetric, so the limiting eigenvectors *are* orthogonal

Basic Concepts

Singular Value Decomposition

We may also take the Singular Value Decomposition (SVD) of A ;

In Matlab : $[U, W, V] = \text{svd}(A)$;

Basic Concepts

Singular Value Decomposition

We may also take the Singular Value Decomposition (SVD) of A ;

In Matlab : $[U, W, V] = \text{svd}(A)$;

It turns out that, we can write the SVD explicitly

$$\begin{aligned} \mathbf{u}_1 &= \frac{1}{L_+} \begin{pmatrix} \epsilon + \sqrt{1 + \epsilon^2} \\ 1 \end{pmatrix}, & \mathbf{u}_2 &= \frac{1}{L_-} \begin{pmatrix} \epsilon - \sqrt{1 + \epsilon^2} \\ 1 \end{pmatrix} \\ w_1 &= \left(2 + \epsilon^2 + 2\sqrt{1 + \epsilon^2}\right)^{1/2}, & w_2 &= \left(2 + \epsilon^2 - 2\sqrt{1 + \epsilon^2}\right)^{1/2} \\ \mathbf{v}_1 &= \frac{1}{L_-} \begin{pmatrix} \sqrt{1 + \epsilon^2} - \epsilon \\ 1 \end{pmatrix}, & \mathbf{v}_2 &= \frac{1}{L_+} \begin{pmatrix} -\sqrt{1 + \epsilon^2} - \epsilon \\ 1 \end{pmatrix} \end{aligned}$$

where $L_+ = \left(1 + (\epsilon + \sqrt{1 + \epsilon^2})^2\right)^{1/2}$, and $L_- = \left(1 + (\epsilon - \sqrt{1 + \epsilon^2})^2\right)^{1/2}$ are normalisations ;

Basic Concepts

Singular Value Decomposition

We may also take the Singular Value Decomposition (SVD) of A ;

In Matlab : $[U, W, V] = \text{svd}(A)$;

It turns out that, we can write the SVD explicitly

$$\begin{aligned} \mathbf{u}_1 &= \frac{1}{L_+} \begin{pmatrix} \epsilon + \sqrt{1 + \epsilon^2} \\ 1 \end{pmatrix}, & \mathbf{u}_2 &= \frac{1}{L_-} \begin{pmatrix} \epsilon - \sqrt{1 + \epsilon^2} \\ 1 \end{pmatrix} \\ w_1 &= \left(2 + \epsilon^2 + 2\sqrt{1 + \epsilon^2}\right)^{1/2}, & w_2 &= \left(2 + \epsilon^2 - 2\sqrt{1 + \epsilon^2}\right)^{1/2} \\ \mathbf{v}_1 &= \frac{1}{L_-} \begin{pmatrix} \sqrt{1 + \epsilon^2} - \epsilon \\ 1 \end{pmatrix}, & \mathbf{v}_2 &= \frac{1}{L_+} \begin{pmatrix} -\sqrt{1 + \epsilon^2} - \epsilon \\ 1 \end{pmatrix} \end{aligned}$$

where $L_+ = \left(1 + (\epsilon + \sqrt{1 + \epsilon^2})^2\right)^{1/2}$, and $L_- = \left(1 + (\epsilon - \sqrt{1 + \epsilon^2})^2\right)^{1/2}$ are normalisations ;

The SVD provides an alternative way to express the matrix inverse, and works also for non square matrices.

Basic Concepts

Properties of SVD for matrices

The SVD of any $N \times M$ matrix A can be written as a product

$$\underbrace{A}_{N \times M} = \underbrace{U}_{N \times N} \underbrace{W}_{N \times M} \underbrace{V^T}_{M \times M} \quad (2.13)$$

The columns of U are orthonormal, as are the columns of V . I.e.

$$U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_N] = \left[\begin{pmatrix} u_{1,1} \\ u_{1,2} \\ \vdots \\ u_{1,N} \end{pmatrix} \begin{pmatrix} u_{2,1} \\ u_{2,2} \\ \vdots \\ u_{2,N} \end{pmatrix} \dots \begin{pmatrix} u_{N,1} \\ u_{N,2} \\ \vdots \\ u_{N,N} \end{pmatrix} \right]$$

and

$$V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_M] = \left[\begin{pmatrix} v_{1,1} \\ v_{1,2} \\ \vdots \\ v_{1,M} \end{pmatrix} \begin{pmatrix} v_{2,1} \\ v_{2,2} \\ \vdots \\ v_{2,M} \end{pmatrix} \dots \begin{pmatrix} v_{M,1} \\ v_{M,2} \\ \vdots \\ v_{M,M} \end{pmatrix} \right]$$

Then we have

$$\mathbf{u}_i \cdot \mathbf{u}_j = \delta_{ij} \quad i, j \in [1 \dots N], \quad \mathbf{v}_i \cdot \mathbf{v}_j = \delta_{ij} \quad i, j \in [1 \dots M].$$

Basic Concepts

Properties of SVD for matrices

The matrix W has a structure depending on the relative values of N, M .

- If $N = M$, W is square, with diagonal elements $W_{i,i} = w_i, i = 1 \dots N$
- If $N > M$, W has a $M \times M$ diagonal matrix above $N - M$ rows of zero :

$$W = \begin{pmatrix} \text{diag}(\mathbf{w}) \\ 0 \end{pmatrix} = \begin{pmatrix} w_1 & 0 & 0 & \dots & 0 \\ 0 & w_2 & 0 & \dots & 0 \\ 0 & 0 & w_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & w_M \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}$$

Basic Concepts

Properties of SVD for matrices

- If $N < M$, W has a $N \times N$ diagonal matrix to the left of $M - N$ columns of zero :

$$W = (\text{diag}(\mathbf{w}) \quad 0) = \begin{pmatrix} w_1 & 0 & 0 & \dots & 0 & 0 \dots & 0 \\ 0 & w_2 & 0 & \dots & 0 & 0 \dots & 0 \\ 0 & 0 & w_3 & \dots & 0 & 0 \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & & \\ 0 & 0 & 0 & \dots & w_N & 0 \dots & 0 \end{pmatrix}$$

The number of non-zeros in the diagonal part of W is equal to the *rank* of the matrix A i.e. the number of linearly independent columns (or rows) of the matrix

$$R = \text{rank}(A) \leq \min(N, M)$$

Basic Concepts

Properties of SVD for matrices

For convenience, assume that the linearly independent columns of U , V are arranged as the first R components, and further that the indexing of columns is such that $|w_i| > |w_{i+1}|$. The elements of the diagonal part of W are termed the *spectrum* of A .

The R column-vectors of V corresponding to the non-zero elements of the diagonal of W form an orthonormal basis spanning a R -dimensional vector space called the *row-space*, and the R column-vectors of U corresponding to these elements form an orthonormal basis spanning a R -dimensional vector space called the *range* $\text{range}(A)$ or *column-space*. These basis functions are related through

$$A\mathbf{v}_i = w_i\mathbf{u}_i, \quad A^T\mathbf{u}_i = w_i\mathbf{v}_i \quad i = 1 \dots R \quad (2.14)$$

The remaining $M - R$ column vectors of V (if greater than zero) form an orthonormal basis spanning a $M - R$ -dimensional vector space called the *null-space* $\text{Null}(A)$, or *kernel* $\text{kern}(A)$. For any vector $\mathbf{x}_\perp \in \text{Null}(A)$ we have

$$A\mathbf{x}_\perp = \mathbf{0} \quad (2.15)$$

Basic Concepts

Properties of SVD for matrices

Similarly, the $N - R$ column vectors of U (if greater than zero) form an orthonormal basis spanning a $N - R$ -dimensional vector space called the *range complement* $\text{range}_{\perp}(A)$ or *co-kernel*. For any vector $\mathbf{b}_{\perp} \in \text{range}_{\perp}(A)$ we have that the linear equation

$$A\mathbf{x} = \mathbf{b}_{\perp} \quad (2.16)$$

has no solutions.

Basic Concepts

Properties of SVD for matrices

For the simple case $N = M = R$ the matrix A is square and full-rank and its inverse can be expressed straightforwardly

$$A^{-1} = VW^{-1}U^T \quad (2.17)$$

where W^{-1} is diagonal with components $1/w_i, i = 1 \dots R$. This representation allows to express the solution of a particular problem

$$A\mathbf{x} = \mathbf{b} \rightarrow \mathbf{x} = \sum_{i=1}^R \frac{\mathbf{u}_i \cdot \mathbf{b}}{w_i} \mathbf{v}_i \quad (2.18)$$

Now we can see a problem that is related to the third definition of an ill-posed problem given above : even if A is full rank, if the relative magnitudes of the w_i decrease rapidly, then the corresponding components of the inverse increase at the same rate. In particular, if the decay of the spectrum is bounded by a geometric ratio then the inverse problem is termed *exponentially ill-posed*.

Basic Concepts

Pseudo-Inverse using SVD

The inverse matrix built from the row-space of A is called the *Moore-Penrose pseudo-inverse* of A .

$$A^\dagger = VW^\dagger U^T = \sum_{i=1}^R \frac{\mathbf{v}_i \mathbf{u}_i^T}{w_i} \quad (2.19)$$

Here W^\dagger is a $M \times N$ matrix with the complementary structure to W . To be specific :

- If $N = M = R$, W^\dagger is square, with diagonal elements $W_{i,i}^\dagger = 1/w_i, i = 1 \dots N$
- If $N > M$, W^\dagger has a $M \times M$ diagonal matrix to the left of $N - M$ columns of zero :

$$W^\dagger = (\text{diag}(\mathbf{1}/\mathbf{w}) \quad \mathbf{0}) = \begin{pmatrix} 1/w_1 & 0 & 0 & \dots & 0 & 0 \dots & 0 \\ 0 & 1/w_2 & 0 & \dots & 0 & 0 \dots & 0 \\ 0 & 0 & 1/w_3 & \dots & 0 & 0 \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & & \\ 0 & 0 & 0 & \dots & 1/w_M & 0 \dots & 0 \end{pmatrix}$$

Basic Concepts

Pseudo-Inverse using SVD

- If $N < M$, W^\dagger has a $N \times N$ diagonal matrix above $M - N$ rows of zero :

$$W^\dagger = \begin{pmatrix} \text{diag}(\mathbf{1}/\mathbf{w}) \\ 0 \end{pmatrix} = \begin{pmatrix} 1/w_1 & 0 & 0 & \dots & 0 \\ 0 & 1/w_2 & 0 & \dots & 0 \\ 0 & 0 & 1/w_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1/w_N \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix}$$

Exercise

Write down the structures of the matrices WW^\dagger and $W^\dagger W$ for the cases $N > M$ and $M < N$ and $M = N$. Assume for the sake of illustration that the rank $R < \min(M, N)$, i.e. that there is a non-trivial null space and range complement space even in the case where $M = N$.

Basic Concepts

Pseudo-Inverse using SVD

The pseudo-inverse can always be constructed regardless if A is non-square and/or rank-deficient. Thus any linear matrix problem has a pseudo-solution

$$A \mathbf{x} = \mathbf{b} \rightarrow \mathbf{x}^\dagger = A^\dagger \mathbf{b} \quad (2.20)$$

This solution has the property that it is the *least-squares, minimum norm* solution in the sense that

$$|\mathbf{b} - A \mathbf{x}^\dagger|^2 \rightarrow \min, |\mathbf{x}^\dagger|^2 \rightarrow \min$$

Basic Concepts

Least Squares-Minimum-Norm property of Pseudo-Inverse using SVD

first we note :

$$\begin{aligned}\mathbf{b} - \mathbf{A}\mathbf{x}^\dagger &= [\mathbf{I} - \mathbf{A}\mathbf{A}^\dagger] \mathbf{b} \\ &= \mathbf{U} [\mathbf{I} - \mathbf{W}\mathbf{W}^\dagger] \mathbf{U}^T \mathbf{b} \\ &= \mathbf{P}_\perp \mathbf{b}\end{aligned}$$

where

$$\mathbf{P}_\perp = \sum_{i=R+1}^N \mathbf{u}_i \mathbf{u}_i^T \quad (2.21)$$

is a matrix that projects data space vectors into the orthogonal complement of the range. Now consider another solution space vector $\tilde{\mathbf{x}} = \mathbf{x}^\dagger + \mathbf{x}'$, with $\mathbf{b}' = \mathbf{A}\mathbf{x}'$ the corresponding data space vector, which by definition lies in the range of \mathbf{A} . We have

$$\begin{aligned}|\mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}|^2 &= |\mathbf{b} - \mathbf{A}\mathbf{x}^\dagger - \mathbf{A}\mathbf{x}'|^2 \\ &= |\mathbf{P}_\perp \mathbf{b} - \mathbf{U}\mathbf{U}^T \mathbf{b}'|^2 \\ &= \sum_{i=R+1}^N (\mathbf{u}_i \cdot \mathbf{b})^2 + \sum_{i=1}^R (\mathbf{u}_i \cdot \mathbf{b}')^2\end{aligned}$$

Basic Concepts

Least Squares-Minimum-Norm property of Pseudo-Inverse using SVD

where the second sum on the right holds due to the fact that $\mathbf{b}' \in \text{Range}(\mathbf{A})$. Thus any vector \mathbf{x}' with a non-zero component in the row space of \mathbf{A} will increase the data space norm.

We now need to see if any vector $\mathbf{x}' \in \text{Null}(\mathbf{A})$ could give a smaller norm $|\tilde{\mathbf{x}}|$. We write

$$\begin{aligned} |\tilde{\mathbf{x}}|^2 &= |\mathbf{x}^\dagger + \mathbf{x}'|^2 \\ &= |\mathbf{V}\mathbf{W}^\dagger\mathbf{U}^T\mathbf{b} + \mathbf{V}\mathbf{V}^T\mathbf{x}'|^2 \\ &= |\mathbf{V}[\mathbf{W}^\dagger\mathbf{U}^T\mathbf{b} + \mathbf{V}^T\mathbf{x}']|^2 \\ &= |[\mathbf{W}^\dagger\mathbf{U}^T\mathbf{b} + \mathbf{V}^T\mathbf{x}']|^2 \\ &= \sum_{i=1}^R \left(\frac{\mathbf{u}_i \cdot \mathbf{b}}{w_i} \right)^2 + \sum_{i=R+1}^M (\mathbf{v}_i \cdot \mathbf{x}')^2 \end{aligned}$$

where the second sum on the right holds due to the fact that $\mathbf{x}' \in \text{Null}(\mathbf{A})$. Thus any vector \mathbf{x}' with a non-zero component in the null-space of \mathbf{A} will increase the solution space norm.

Basic Concepts

Analysing Problem 1 using SVD

Problem 1 revisited

$$A = \begin{pmatrix} 1 & 1 \end{pmatrix} \Rightarrow A^\dagger = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\text{SVD:} \quad (1) (\sqrt{2} \ 0) \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = A^\dagger(2) = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

The row-space is spanned by vector $\frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, and the null space is spanned by vector $\frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$.

Basic Concepts

Analysing Problem 2 using SVD

Problem 2 revisited

$$A = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow A^\dagger = \frac{1}{2} \begin{pmatrix} 1 & 1 \end{pmatrix}$$

$$\text{SVD:} \quad \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \sqrt{2} \\ 0 \end{pmatrix} (1)$$

$$\Rightarrow x_1 = A^\dagger \begin{pmatrix} 1 \\ 3 \end{pmatrix} = 2$$

The range is spanned by vector $\frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, and the complement of the range is spanned by vector $\frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$.

Basic Concepts

Analysing Problem 3 using SVD

Problem 3 revisited

The matrix inverse was given in eq. 2.12 so we do not need to write it again. However, by considering it as

$$A^\dagger(\mathbf{b}_0 + \boldsymbol{\eta}) = A^\dagger \mathbf{b}_0 + \frac{\boldsymbol{\eta} \cdot \mathbf{u}_1}{w_1} \mathbf{v}_1 + \frac{\boldsymbol{\eta} \cdot \mathbf{u}_2}{w_2} \mathbf{v}_2$$

We see that the solution can be considered as a random vector with mean the noise free vector $A^\dagger \mathbf{b}_0$ and covariance given by

$$\mathbf{C} = \frac{\mathbf{v}_1 \mathbf{v}_1^T}{w_1^2} + \frac{\mathbf{v}_2 \mathbf{v}_2^T}{w_2^2} = \mathbf{V} \mathbf{W}^{-2} \mathbf{V}^T$$

Basic Concepts

Singular Value Decomposition for Operators

Most of what we have discussed for matrices applies also to Operators. Let us restrict the discussion to linear integral operators $A : X \mapsto Y$

$$b = Af \quad \equiv \quad b(y) = \int_{-\infty}^{\infty} K(y, x) f(x) dx \quad (2.22)$$

We define the *adjoint* operator $A^* : Y \mapsto X$ (the equivalent of matrix transpose) by integration of the (complex conjugate of the) kernel function $K(x, y)$ with the other variable :

$$h = A^* g \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} \bar{K}(y, x) g(y) dy \quad (2.23)$$

and we also need to attach a meaning to the inner product :

$$\langle h, f \rangle_X := \int_{-\infty}^{\infty} \bar{h}(x) f(x) dx, \quad \langle g, b \rangle_Y := \int_{-\infty}^{\infty} \bar{g}(y) b(y) dy \quad (2.24)$$

Note that the inner product is technically different for the functions in the data space and the solution space, because they might in general have different variables.

Basic Concepts

Singular Value Decomposition for Operators

The SVD of operator A is then defined as the set of *singular functions* $\{v_k(x), u_k(y)\}$ and singular values $\{w_k\}$ such that

$$Av_k = w_k u_k; \quad A^* u_k = w_k v_k$$

and with

$$\langle u_j, u_k \rangle_Y = \delta_{jk}, \quad \langle v_j, v_k \rangle_X = \delta_{jk}.$$

One major difference is that the set of such functions is (usually) *infinite*. There is often no equivalence to the kernel (Null space) of the operator, although sometimes there is.

Basic Concepts

Singular Value Decomposition for Operators

Note also the mappings

$A^*A : X \mapsto X$ given by

$$h = A^*Af \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \bar{K}(y, x) K(y, x') f(x') dx' dy$$

$$\equiv \quad h(x) = \int_{-\infty}^{\infty} B(x, x') f(x') dx'$$

$$\text{where} \quad B(x, x') = \int_{-\infty}^{\infty} \bar{K}(y, x) K(y, x') dy$$

$AA^* : Y \mapsto Y$ given by

$$b = AA^*g \quad \equiv \quad b(y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \bar{K}(y', x) K(y, x) g(y') dx dy$$

$$\equiv \quad b(y) = \int_{-\infty}^{\infty} C(y, y') g(y') dy'$$

$$\text{where} \quad C(y, y') = \int_{-\infty}^{\infty} \bar{K}(y', x) K(y, x) dx$$

Basic Concepts

Example Singular Value Decomposition for Operators

Example: convolution as a linear operator with stationary kernel. By *stationary* we mean that the kernel depends only on the difference $(x - y)$, in other words $K(x, y)$ is a function of only one variable $K(y - x) \equiv K(t)$

$$b = Af := K * f \quad \equiv \quad b(y) = \int_{-\infty}^{\infty} K(y - x)f(x)dx \quad (2.25)$$

An important property of convolution is the *Convolution Theorem* which relates the functions under convolution, and the result, to their Fourier Transforms. Assuming that

$$\hat{b}(k) = \mathcal{F}_{x \rightarrow k} b(x), \quad \hat{K}(k) = \mathcal{F}_{x \rightarrow k} K(x), \quad \hat{f}(k) = \mathcal{F}_{x \rightarrow k} f(x)$$

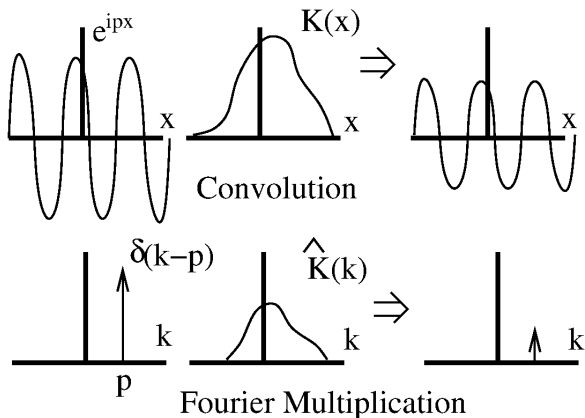
this convolution theorem states that

$$b = K * f \quad \Rightarrow \quad \hat{b} = \hat{K}\hat{f} \quad (2.26)$$

Basic Concepts

Example Singular Value Decomposition for Operators

Convolution of a trigonometric function of frequency p with any function K results in the same frequency trigonometric function, scaled by amplitude $\hat{K}(p)$



Basic Concepts

Example Singular Value Decomposition for Operators

It follows from the convolution theorem in Fourier theory that the singular vectors of A are the functions $\{\cos px, \sin px\}$ with singular values $\hat{K}(p)$ where \hat{K} is the Fourier Transform of K

$$A \begin{pmatrix} \cos px \\ \sin px \end{pmatrix} = \hat{K}(p) \begin{pmatrix} \cos py \\ \sin py \end{pmatrix} \quad (2.27)$$

$$A^* \begin{pmatrix} \cos py \\ \sin py \end{pmatrix} = \hat{K}(p) \begin{pmatrix} \cos px \\ \sin px \end{pmatrix} \quad (2.28)$$

Exercise

Verify that $\sin px$ and $\cos px$ are eigenvectors (and therefore also singular vectors) of the linear convolution operator with a Gaussian

$K(y - x) = \exp \left[-\frac{(y-x)^2}{2\sigma^2} \right]$ (for any non zero value of σ) and that singular values are $\lambda_p = \exp \left[-\frac{1}{2}\sigma^2 p^2 \right]$.

Basic Concepts

Example Singular Value Decomposition for Operators

We also see exactly how to do the inversion - applying the Moore-Penrose inverse in the continuous case is given by

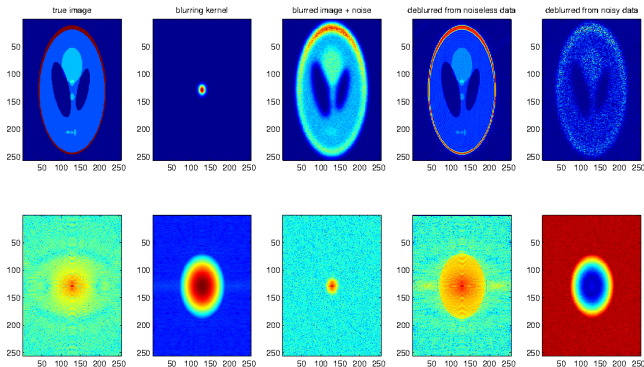
$$f^{\text{recon}}(x) = \sum_{i=1}^{\infty} v_i(x) \frac{\langle b(x), u_i(x) \rangle}{w_i} \equiv \mathcal{F}_{k \rightarrow x}^{-1} \left[\frac{\mathcal{F}_{x \rightarrow k} b(x)}{\mathcal{F}_{x \rightarrow k} K(x)} \right]$$

This also makes clear how the noise is propagated - the covariance of noise in the reconstruction is diagonalised by its Fourier components, and the higher spatial frequencies are amplified by the reciprocal of the square of the corresponding frequency of the convolution kernel. An example is shown in the next slide. The convolution kernel is a Gaussian so the noise is propagated with the exponential of the square of the frequency number.

Basic Concepts

Example Singular Value Decomposition for Operators

Convolution and Deconvolution in the Fourier Domain. Top row is spatial domain and bottom row is the Fourier Transform of the top row, on a logarithmic scale.



Basic Concepts

The Continuous-Discrete Case

Lastly, we may consider the case where the data are a finite set of samples, but the reconstruction space consists of functions. The forward mapping is

$$\mathbf{b} = A\mathbf{f} \quad \equiv \quad \mathbf{b} = \int_{-\infty}^{\infty} \mathbf{K}(x)f(x)dx$$

where $\mathbf{K}(x)$ is a vector valued function. The Adjoint becomes a sum of functions

$$h = A^*\mathbf{g} \quad \equiv \quad h(x) = \mathbf{g}^T \mathbf{K}(x) = \sum_{i=1}^N g_i \bar{K}_i(x)$$

Basic Concepts

The Continuous-Discrete Case

We also see that $A^*A : X \mapsto X$ is an operator

$$h = A^*A f \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} \underbrace{\mathbf{K}^T(x) \mathbf{K}(x')}_{B(x, x')} f(x') dx' = \int_{-\infty}^{\infty} \underbrace{\left(\sum_{i=1}^N \bar{K}_i(x) K_i(x') \right)}_{B(x, x')} f(x') dx'$$

and that $AA^* : Y \mapsto Y$ is a $N \times N$ matrix

$$\mathbf{b} = AA^* \mathbf{g} \quad \equiv \quad \mathbf{b} = \underbrace{\left(\int_{-\infty}^{\infty} \mathbf{K}(x) \mathbf{K}^T(x) dx \right)}_{\mathbf{C}} \mathbf{g} = \sum_{j=1}^N \underbrace{\left(\int_{-\infty}^{\infty} K_i(x) K_j^T(x) dx \right)}_{C_{ij}} g_j$$

In this case the number of left singular vectors $\{\mathbf{u}_i; i = 1 \dots N\}$ is finite, and the number of right singular functions $\{v_i(x)\}$ is (usually) infinite, but those greater than $i = N$ will be in the null-space of the forward operator. This is the *usual* case for real measurements : there exists an infinite dimensional space Null-Space, “invisible” to the measurements!