

# Optimisation and Inverse Problems in Imaging

Simon R. Arridge

February 4, 2017

† *Department of Computer Science, University College London, Gower Street London, WC1E 6BT*

**GV08**

*2007-2016*

# 1 Introduction

The topic of Optimisation covers many areas. In this course we are going to consider the role of Optimisation in imaging. In particular (but not exclusively) we are going to consider the role of optimisation in the topic of *Inverse Problems*, and within this subject we are going to consider in particular (but not exclusively) the topic of *Inverse Problems in Imaging*.

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?”

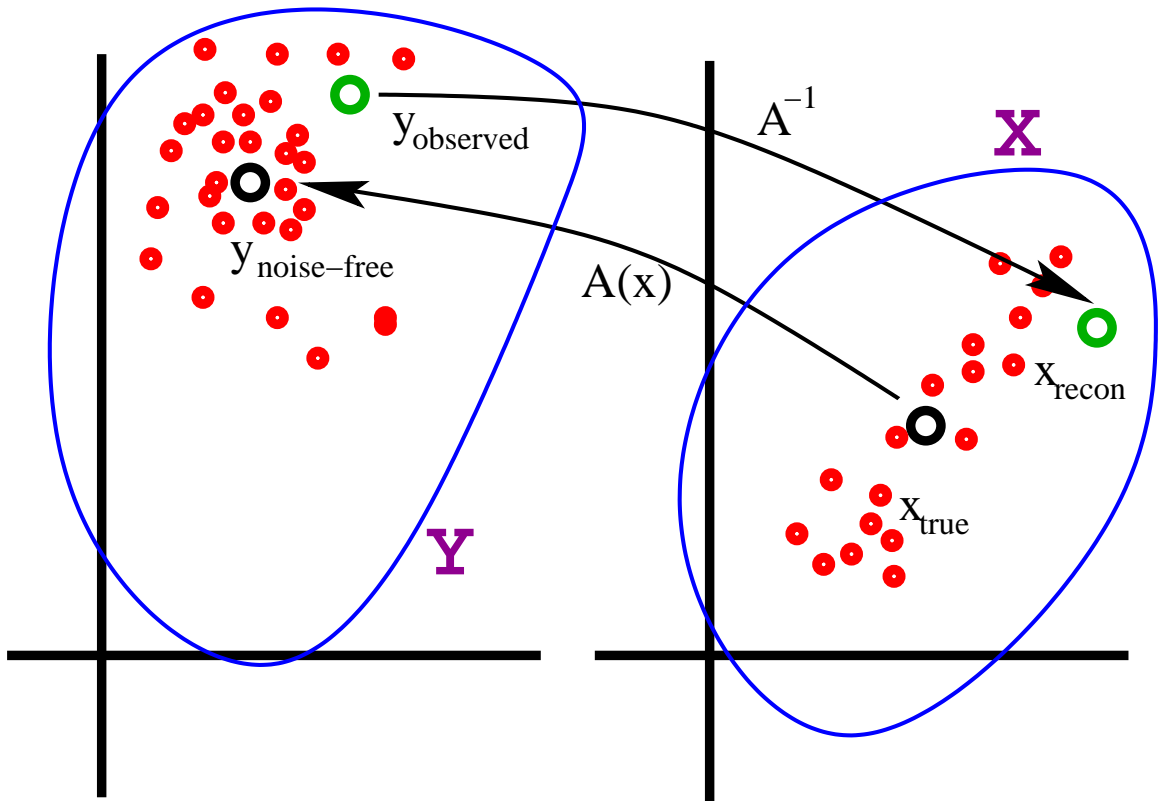


Figure 1.1:

There are many generalisations. We may sometimes consider not one but multiple metrics. An important type of optimisation is *model-fitting*. In this case we need a 5th concept: *Data*. In this case the metric is not just applied to the model, but to the model and the data together to evaluate how well the model predicts the data.

## 2 Basic Concepts

### 2.1 Some examples

#### 2.1.1 Minimum norm solutions to Underdetermined Problems

Consider the following problem : Solve

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

There is no unique solution. The set of possible solutions forms a line on the  $x_1, x_2$  plane. However the following can be solved :

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

which has solution  $(x_1, x_2) = (1, 2)$ . It is illustrated graphically in figure 2.1. This is an example

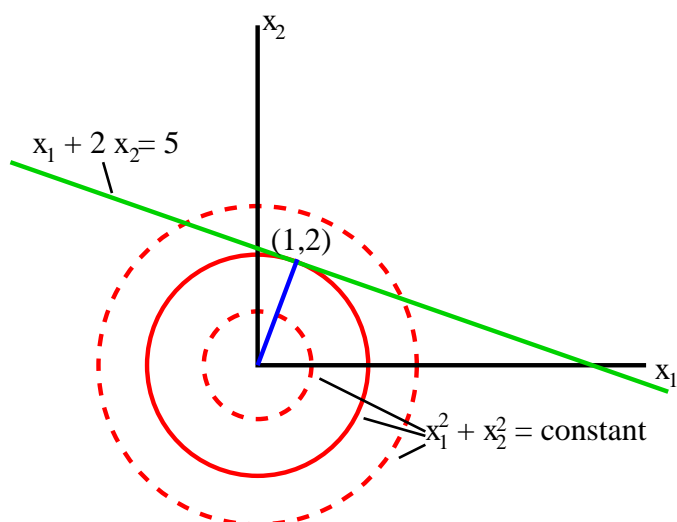


Figure 2.1:

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  $\hat{\mathbf{x}}_1 = (1, 0)$ ,  $\hat{\mathbf{x}}_2 = (0, 1)$ . In this

case, because the norm is quadratic, the solution is found by taking the *Moore-Penrose generalised inverse* of the underdetermined linear equation Eq. 2.1 :

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \left( (12) \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right)^{-1} (5) = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (2.3)$$

and the solution is the *minimum norm* solution, where the norm is the normal Euclidean  $L_2$  norm.

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

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

When  $1 < p < \infty$  the surfaces  $\Phi = \text{constant}$  are *convex* and solutions can be found relatively easily. 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  $p = \infty$ .

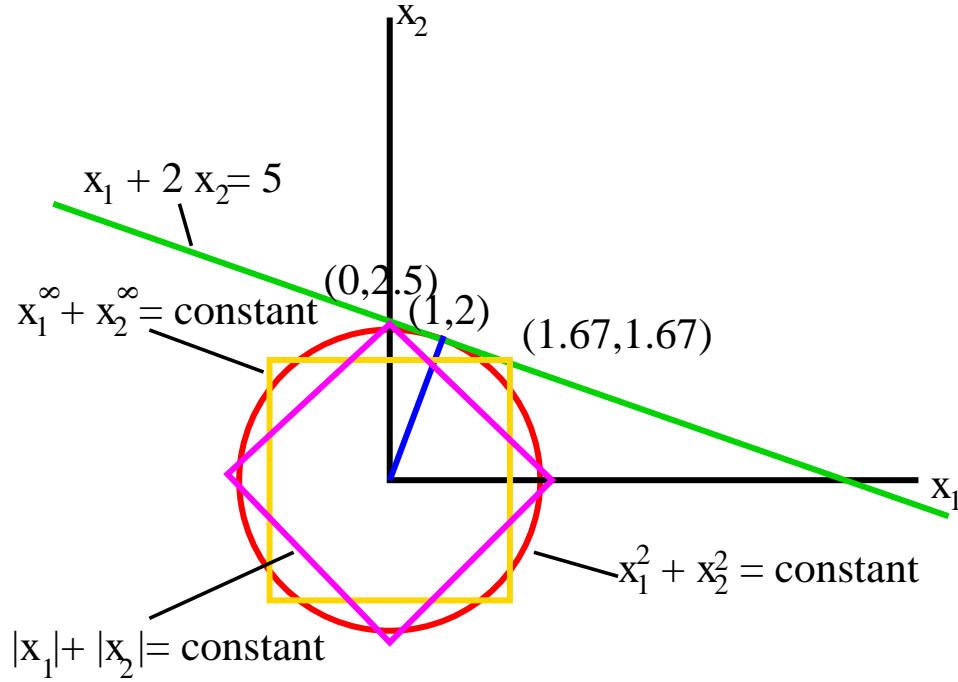


Figure 2.2:

However, actually calculating the solution is not as simple as solving a linear problem like Eq. 2.3. The case  $p = 1$  is known as the *sparse* solution.

The case  $0 < p < 1$  leads to *concave* surfaces and to the more complex problem of *non-convex* optimisation.

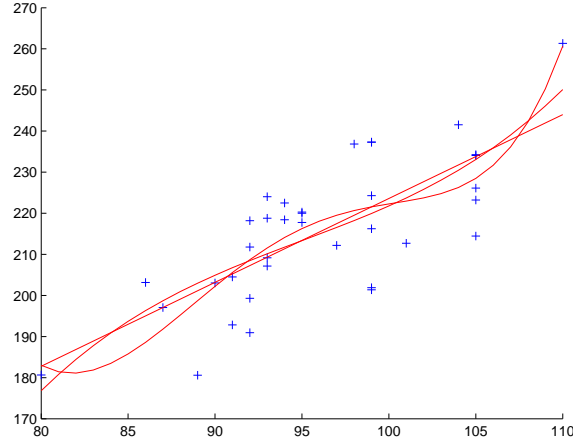


Figure 2.3: Fitting polynomials of increasing order to a set of points

### 2.1.2 Model Fitting

Suppose we have a set of data  $\{x_i, y_i; i = 1 \dots N\}$  and a hypothetical model

$$y = A(x; \boldsymbol{\xi})$$

parameterised by values  $\{\xi_0, \xi_1, \dots, \xi_M\}$ . We can provide a fit of the model by optimising a data fitting function such as

$$\Phi = \frac{1}{2} \sum_{i=1}^N |y_i - A(x_i, \boldsymbol{\xi})|^2 \quad (2.5)$$

see figure 2.3. More generally, we may define a *data-fitting* functional

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

This framework admits, for example, the  $p$ -norms we introduced in Eq.(2.4), or they may be more general norms such as Kullback-Leibler divergence, or Cross-entropy. In general the choice of norm in the space of the *data* is dictated by the expected statistics of the errors in the data. A choice of the  $L_2$  norm corresponds to assuming *identically distributed, independent Gaussian random variable* for the errors (iid Normal distributions).

If the model is linear, and the data-fitting functional is sum of squares the problem is *linear least-squares* and can be solved using linear algebra. For example, fitting a cubic polynomial  $y = \xi_0 + \xi_1 x + \xi_2 x^2 + \xi_3 x^3$  to  $N > 4$  data is an *over-determined* problem defined by

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{pmatrix} \begin{pmatrix} \xi_0 \\ \xi_1 \\ \xi_2 \\ \xi_3 \end{pmatrix} \quad (2.7)$$

If  $A$  is non-linear with respect to the parameters  $\xi$  then the process is one of *non-linear, unconstrained optimisation*. This could be non-linear least-squares, or a more general problem.

### 2.1.3 Inequality constraints

Consider the same problem as in section 2.1.1, but introduce bound on the solution

$$\begin{array}{ll} \text{minimise} & \Phi = x_1^2 + x_2^2 \\ \text{subject to} & x_1 + 2x_2 = 5 \end{array} \quad (2.8)$$

$$\begin{array}{l} 2 \leq x_1 \leq 4 \\ 1 \leq x_2 \leq 5 \end{array} \quad (2.9)$$

which has solution  $(x_1, x_2) = (2, 1.5)$ , see figure 2.4. Unlike Eq. 2.2 this cannot be solved by a

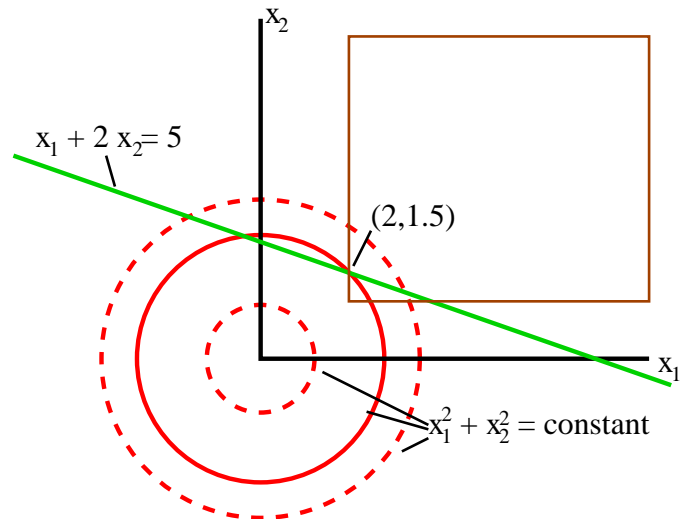


Figure 2.4:

simple linear method like Eq. 2.3.

#### References for this section

[1, 2, 3, 4]

## 2.2 Inverse Problems

### 2.2.1 Problem 1 : Non Uniqueness

Solve

$$x_1 + x_2 = 2 \quad (2.10)$$

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.

### 2.2.2 Problem 2 : Non Existence

Solve

$$\begin{aligned}x_1 &= 1 \\x_1 &= 3\end{aligned}\tag{2.11}$$

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

$$\Phi = (x_1 - 1)^2 + (x_1 - 3)^2 \rightarrow \min\tag{2.12}$$

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. Instead we can find the minimum :  $x_1 = 2, \Phi = 2$ . This is the *least squares solution*.

### 2.2.3 Problem 3 : Discontinuous dependence on data

Solve

$$A\mathbf{x} = \begin{pmatrix} 1 & 1 + \epsilon \\ 1 - \epsilon & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}\tag{2.13}$$

This problem has a solution that tends towards instability. There are a number of ways to see this. The matrix  $A$  has determinant  $\epsilon^2$  so its inverse is

$$A^{-1} = \frac{1}{\epsilon^2} \begin{pmatrix} 1 & -(1 + \epsilon) \\ -(1 - \epsilon) & 1 \end{pmatrix}\tag{2.14}$$

Thus the inverse grows as an inverse square in the parameter  $\epsilon$ .

To see the effect of the inverse with respect to noise, consider the following example :

Results are shown in figure 2.2.3. We see that the data has an isotropic distribution around the specified noise free point whereas the solution has a highly *non-isotropic* distribution. 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  $Ax_i = b_i$

**end for**

---

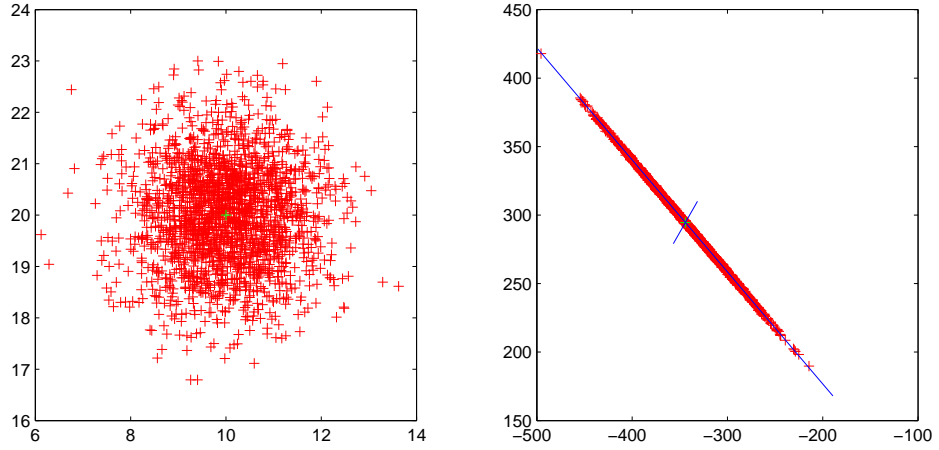


Figure 2.5: Results of ill-posed matrix inversion trial for  $\epsilon \simeq 0.2$ . Left : The distribution of data points. Right : the distribution of solutions. The singular values of the matrix are 2.02, 0.0202, and the standard deviation of the covariance of the solution points are 1.9984, 0.020502

highly anisotropic amplification of noise preferentially in some directions is an example of the third kind of ill-posedness of inverse problems. To understand the nature of this ill-posedness we can investigate the structure of the linear forward problem using the Singular Value Decomposition (SVD).

### References for this section

[5]

## 2.3 Singular Value Decomposition (SVD)

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)

$$e_1 = \begin{pmatrix} \sqrt{1 - \epsilon^2} \\ 1 - \epsilon \end{pmatrix}, \quad e_2 = \begin{pmatrix} -\sqrt{1 - \epsilon^2} \\ 1 - \epsilon \end{pmatrix}$$



[NOTE : Since the matrix 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 the matrix is symmetric, so the limiting eigenvectors are orthogonal]

We may also take the Singular Value Decomposition (SVD); In Matlab :

```
eps = 0.01;
A = [1 1+eps; 1 - eps 1];
[U,W,V] = svd(A);
```

[Note that Matlab requires numerics. The symbolic part of Matlab is not very comprehensive]

We can write the SVD explicitly

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

where  $L_+ = (1 + (\epsilon + \sqrt{1 + \epsilon^2})^2)^{1/2}$ , and  $L_- = (1 + (\epsilon - \sqrt{1 + \epsilon^2})^2)^{1/2}$  are normalisations ; [Note that normalisation must be carried out for SVD to provide the correct mappings between left and right spaces].

**Exercise 1** Verify the above SVD. I.e. show

$$\mathbf{A}\mathbf{v}_i = w_i\mathbf{u}_i, \quad \mathbf{A}^T\mathbf{u}_i = w_i\mathbf{v}_i \quad i = 1, 2$$

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

### 2.3.1 Properties of SVD for matrices

The SVD can be defined for any linear operator. To begin with we will consider matrices. The SVD of any  $N \times M$  matrix  $\mathbf{A}$  can be written as a product

$$\underbrace{\mathbf{A}}_{N \times M} = \underbrace{\mathbf{U}}_{N \times N} \underbrace{\mathbf{W}}_{N \times M} \underbrace{\mathbf{V}^T}_{M \times M} \quad (2.15)$$

The columns of  $\mathbf{U}$  are orthonormal, as are the columns of  $\mathbf{V}$ . I.e. if we consider

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

and

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

Then we have

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

The matrix  $\mathbf{W}$  has a structure depending on whether  $N$  is greater than, equal to, or less than  $M$ .

If  $N = M$ ,  $\mathbf{W}$  is square, with diagonal elements  $W_{i,i} = w_i, i = 1 \dots N$

If  $N > M$ ,  $\mathbf{W}$  has a  $M \times M$  diagonal matrix above  $N - M$  rows of zero :

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

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

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

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

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

For convenience, assume that the linearly independent columns of  $\mathbf{U}$ ,  $\mathbf{V}$  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  $\mathbf{W}$  are termed the *spectrum* of  $\mathbf{A}$ .

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

$$\mathbf{A}\mathbf{v}_i = w_i\mathbf{u}_i, \quad \mathbf{A}^T\mathbf{u}_i = w_i\mathbf{v}_i \quad i = 1 \dots R \quad (2.16)$$

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

$$\mathbf{A}\mathbf{x}_\perp = \mathbf{0} \quad (2.17)$$

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

$$\mathbf{A}\mathbf{x} = \mathbf{b}_\perp \quad (2.18)$$

has no solutions.

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

$$\mathbf{A}^{-1} = \mathbf{V}\mathbf{W}^{-1}\mathbf{U}^T \quad (2.19)$$

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

$$\mathbf{A}\mathbf{x} = \mathbf{b} \rightarrow \mathbf{x} = \sum_{i=1}^R \frac{\mathbf{u}_i \cdot \mathbf{b}}{w_i} \mathbf{v}_i \quad (2.20)$$

Now we can see a problem that is related to the third definition of an ill-posed problem given above : even if  $\mathbf{A}$  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*.

### 2.3.2 Pseudo-Inverse

The inverse matrix built from the row-space of  $\mathbf{A}$  is called the *Moore-Penrose pseudo-inverse* of  $\mathbf{A}$ .

$$\mathbf{A}^\dagger = \mathbf{V}\mathbf{W}^\dagger\mathbf{U}^T = \sum_{i=1}^R \frac{\mathbf{v}_i\mathbf{u}_i^T}{w_i} \quad (2.21)$$

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

If  $N = M = R$ ,  $\mathbf{W}^\dagger$  is square, with diagonal elements  $W_{i,i}^\dagger = 1/w_i, i = 1 \dots N$

If  $N > M$ ,  $\mathbf{W}^\dagger$  has a  $M \times M$  diagonal matrix to the left of  $N - M$  columns of zero :

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

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

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

**Exercise 2** Write down the structures of the matrices  $\mathbf{W}\mathbf{W}^\dagger$  and  $\mathbf{W}^\dagger\mathbf{W}$  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  $M = N$ .

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

$$\mathbf{A}\mathbf{x} = \mathbf{b} \rightarrow \mathbf{x}^\dagger = \mathbf{A}^\dagger\mathbf{b} \quad (2.22)$$

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

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

To see this, first we note :

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

where

$$\mathbf{P}_\perp = \sum_{i=R+1}^N \mathbf{u}_i \mathbf{u}_i^\mathbf{T} \quad (2.23)$$

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

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

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 Null-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

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

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.

To complete our exposition, let us solve the three problems in section 2.2 using the pseudo-inverse.

### 2.3.3 Problem 1

$$\begin{aligned}
\mathbf{A} = \begin{pmatrix} 1 & 1 \end{pmatrix} &\Rightarrow \mathbf{A}^\dagger = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\
\text{SVD:} & \quad (1) \begin{pmatrix} \sqrt{2} & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \\
\Rightarrow \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \mathbf{A}^\dagger(2) &= \begin{pmatrix} 1 \\ 1 \end{pmatrix}
\end{aligned}$$

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  $\frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$ .

### 2.3.4 Problem 2

$$\begin{aligned}
\mathbf{A} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} &\Rightarrow \mathbf{A}^\dagger = \frac{1}{2} \begin{pmatrix} 1 & 1 \end{pmatrix} \\
\text{SVD:} & \quad \begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} \sqrt{2} \\ 0 \end{pmatrix} (1) \\
\Rightarrow x_1 = \mathbf{A}^\dagger \begin{pmatrix} 1 \\ 3 \end{pmatrix} &= 2
\end{aligned}$$

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  $\frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$ .

### 2.3.5 Problem 3

The matrix inverse was given in Eq. 2.14 so we do not need to write it again. However, by considering it as

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

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

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

## 2.4 SVD for Operators

Most of what we have discussed for matrices applies also to Linear Operators.

Let us restrict the discussion to linear integral operators  $A : X \mapsto Y$

$$b = A f \quad \equiv \quad b(y) = \int_{-\infty}^{\infty} K(y, x) f(x) dx \quad (2.24)$$

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 :

$$h = A^* g \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} \overline{K}(y, x) g(y) dy \quad (2.25)$$

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

$$\langle h, f \rangle_X := \int_{-\infty}^{\infty} \overline{h}(x) f(x) dx, \quad \langle g, b \rangle_Y := \int_{-\infty}^{\infty} \overline{g}(y) b(y) dy \quad (2.26)$$

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 (in this case they are the same :  $x$  and  $y$ ).

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

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

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. Note also the mappings

$$\begin{aligned} A^* A : X \mapsto X & \quad \text{given by} \\ h = A^* A f & \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \overline{K}(y, x) K(y, x') f(x') dx' dy \\ & \quad \equiv \quad h(x) = \int_{-\infty}^{\infty} B(x, x') f(x') dx' \\ \text{where} \quad B(x, x') & = \int_{-\infty}^{\infty} \overline{K}(y, x) K(y, x') dy \end{aligned}$$

and

$$\begin{aligned}
AA^* : Y &\mapsto Y && \text{given by} \\
b = A^*Ag &\equiv && b(y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \overline{K}(y', x) K(y, x) g(y') dx dy \\
&\equiv && b(y) = \int_{-\infty}^{\infty} C(y, y') g(y') dy' \\
\text{where } C(y, y') &= \int_{-\infty}^{\infty} \overline{K}(y', x) K(y, x) dx
\end{aligned}$$

**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  $K(y - x) \equiv K(t)$

$$b = Af := K * f \quad \equiv \quad b(y) = \int_{-\infty}^{\infty} K(y - x) f(x) dx \quad (2.27)$$

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

$$b = K * f \quad \Rightarrow \quad \hat{b} = \hat{K} \hat{f} \quad (2.28)$$

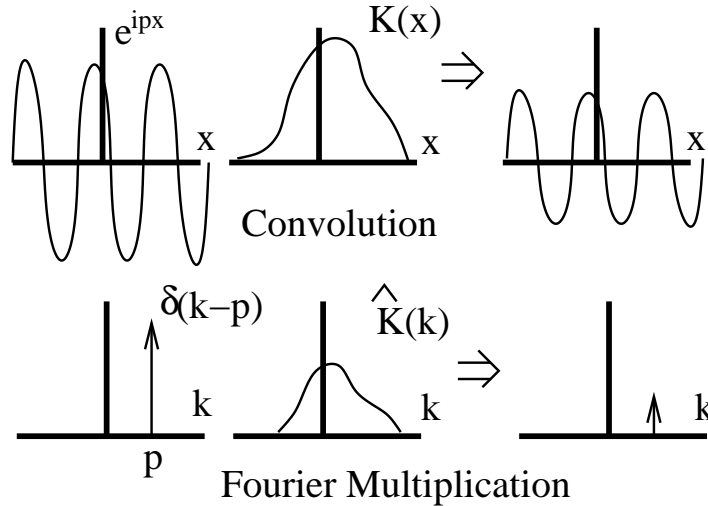


Figure 2.6: Convolution of a trigonometric function of frequency  $p$  with any function  $K$  results in the same frequency trigonometric function, scaled by amplitude  $\hat{K}(p)$

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  $K$

$$A \begin{pmatrix} \cos px \\ \sin px \end{pmatrix} = \hat{K}(p) \begin{pmatrix} \cos py \\ \sin py \end{pmatrix} \quad (2.29)$$

$$A^* \begin{pmatrix} \cos py \\ \sin py \end{pmatrix} = \hat{K}(p) \begin{pmatrix} \cos px \\ \sin px \end{pmatrix} \quad (2.30)$$

**Exercise 3** 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  $\sigma$ ) and that singular values are  $\lambda_p = \exp\left[-\frac{1}{2}\sigma^2 p^2\right]$ .

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 show in figure 2.4; the convolution kernel is a Gaussian so the noise is propagated with the exponential of the square of the frequency number.

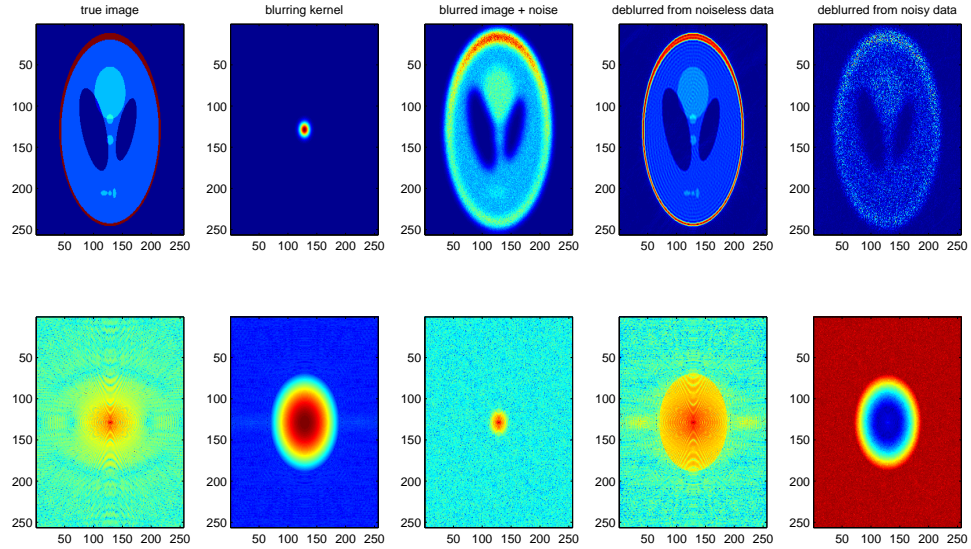


Figure 2.7: Convolution and Deconvolution in the Fourier Domain. Top row is spatial domain and bottom row is the Fourier Transform of the top row, on a logarithmic scale.

## 2.5 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

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

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

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

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!

### References for this section

Chapter 2 of [6], Chapter 4 of [7], Chapter 9 of [8].

### 3 Introduction to Regularisation

#### 3.1 Illustrative Example

Consider the Fredholm integral equation of the first kind representing stationary (space-invariant) convolution in one space dimension

$$g(x) = \mathcal{A}f := \int_0^1 k(x - x')f(x')dx' \quad x \in [0, 1] \quad (3.1)$$

The kernel  $k(x)$  is a blurring function, which is a function of only one parameter because we defined the problem to be space-invariant. In imaging systems it is called a *Point Spread Function* (PSF) [note the connection with the *Impulse Response Function* mentioned in section section B.1.3, Eq. B.3].

For an illustrative example take the kernel to be a Gaussian

$$k(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{x^2}{2\sigma^2}\right] \quad (3.2)$$

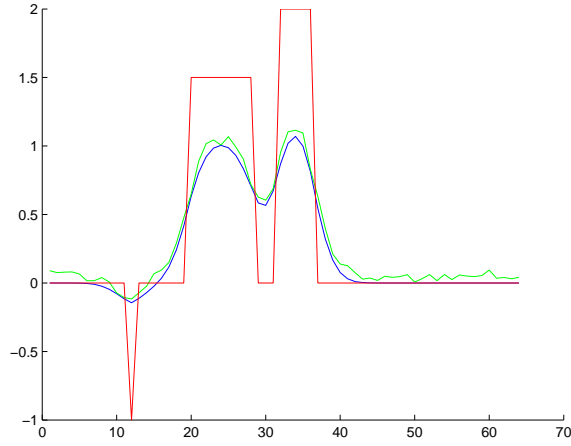


Figure 3.1: Example of blurring by Eq. 3.1 with kernel Eq. 3.2. Red curve : the original function; Blue curve the output  $\mathcal{A}f$ ; green curve the data with added random noise.

The effect of the blurring is shown in figure 3.1. The data is corrupted with some random “measurement” noise to give data

$$\tilde{g}(x) = \mathcal{A}f + \eta \quad (3.3)$$

We want to solve the inverse problem to find  $f$  given  $\tilde{g}$ . To introduce the ideas, we will use a

discrete setting. We construct a discrete matrix approximation  $\mathbf{A}$  to  $\mathcal{A}$  as

$$A_{ij} = \frac{h}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{((i-j)h)^2}{2\sigma^2} \right] \quad i, j \in [1 \dots n] \quad (3.4)$$

where  $h = 1/n$  is the discretisation size.

The matrix  $\mathbf{A}$  is rank-complete so is invertible.

$$f^\dagger = \mathbf{A}^\dagger \tilde{g} = \mathbf{A}^{-1} \tilde{g} \quad (3.5)$$

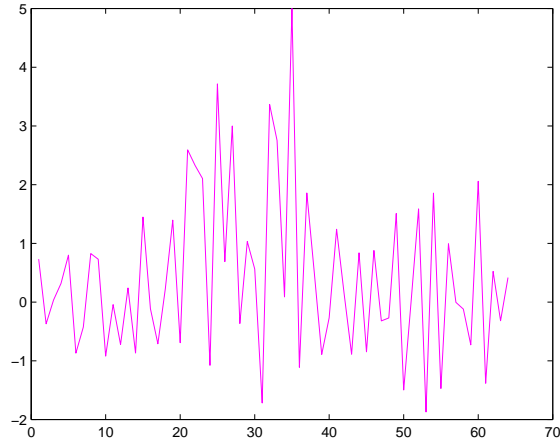


Figure 3.2: Naive reconstruction of  $f$  from  $\tilde{g}$  without regularisation.

Figure 3.2 shows what happens if we apply the Moore-Penrose inverse directly to the noisy data. The reason is seen from the SVD of  $\mathbf{A}$ . The singular spectrum decays exponentially (in fact as  $\exp(-i^2)$ ), so the inverse adds components of the noise whose projections on the singular vectors with higher index increase exponentially.

### 3.2 Regularisation by Filtering

Our first idea for regularisation comes from the idea of *filtering* the contribution of the singular vectors in the inverse. Rather than the pseudo inverse Eq. 2.21 we construct a *regularised inverse*

$$\mathbf{A}_\alpha^\dagger = \mathbf{V} \mathbf{W}_\alpha^\dagger \mathbf{U}^T = \sum_{i=1}^R \mathbf{v}_i \frac{q_\alpha(w_i^2)}{w_i} \mathbf{u}_i^T \quad (3.6)$$

where the function  $q_\alpha(w_i^2)$  is a filter on the SVD spectrum. There are two commonly used filters :

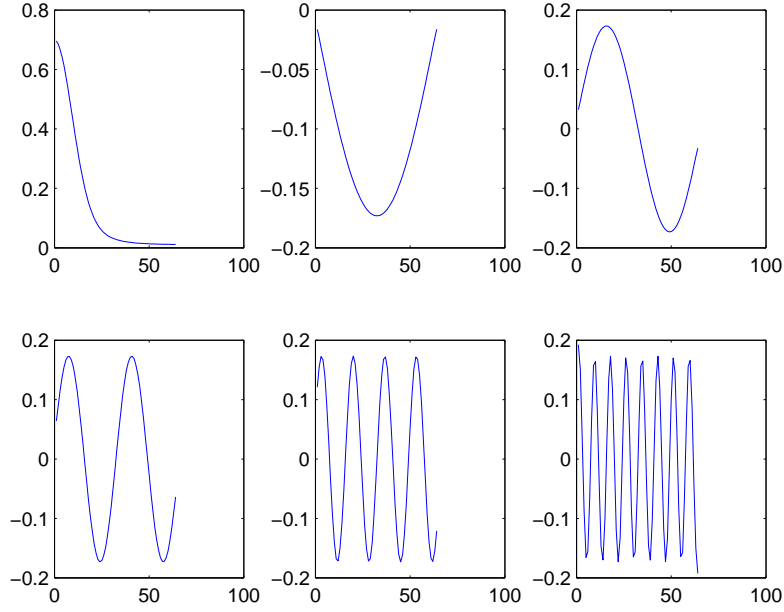


Figure 3.3: SVD of the matrix  $A$ . Shown are the singular spectrum  $\{w_i\}$  and the singular vectors 1,2,4,8,16. The matrix is symmetric so the left and right singular vectors are identical.

**Truncated SVD** : here the filter is

$$q_\alpha(w_i^2) = \begin{cases} 1, & \text{if } i < \alpha \\ 0, & \text{if } i \geq \alpha \end{cases} \quad (3.7)$$

The parameter  $\alpha$  acts as a threshold, cutting off the contribution of higher order singular vectors. This is exactly analogous to the idea of a *low-pass filter* in signal processing, and the analysis of SVD filtering sometimes goes by the name of *Generalised Fourier Analysis*. The regularised inverse in this case is

$$A_\alpha^\dagger = \sum_{i=1}^{\alpha} \frac{\mathbf{v}_i \mathbf{u}_i^T}{w_i} \quad (3.8)$$

Rather than supply  $\alpha$  as an index on the singular values (which requires them to be ordered in monotonically decreasing sequence), we could specify the threshold on the value of the singular component

$$q_\alpha(w_i^2) = \begin{cases} 1, & \text{if } w_i^2 > \alpha \\ 0, & \text{if } w_i^2 \leq \alpha \end{cases} \quad (3.9)$$

**Tikhonov filtering** : here the filter is

$$q_\alpha(w_i^2) = \frac{w_i^2}{w_i^2 + \alpha} \quad (3.10)$$

The regularised inverse in this case is

$$\mathbf{A}_\alpha^\dagger = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \quad (3.11)$$

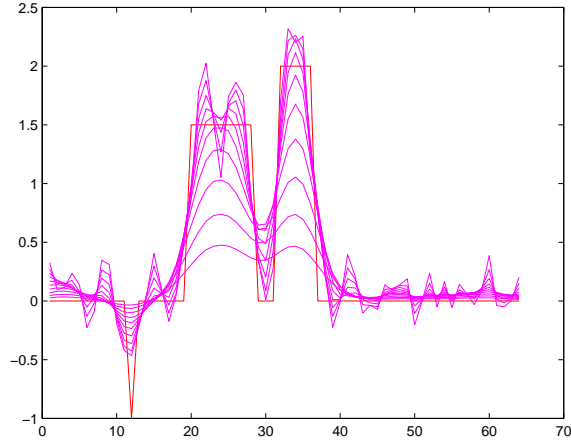


Figure 3.4: Reconstruction of  $f$  from  $\tilde{g}$  with Tikhonov regularisation. Curves show the effect of decreasing  $\alpha$  from 1 in steps of  $\frac{1}{2}$

In figure 3.4 we show the effect of Tikhonov regularisation with decreasing parameter  $\alpha$ . The higher values give a strong smoothing effect; the smaller values recover more of the function, at the expense of increasing noise.

**Exercise 4** The deblurring problem of section 3.1 is very similar but not identical to smoothing with a low-pass filter. Why is it different? Repeat the results of this section using a low-pass filter implemented in the Fourier domain, and comment on the differences.

### 3.3 Variational Regularisation Methods

When the system under consideration is large it is usually impractical to construct the SVD explicitly. However we can specify the solution by minimisation of a variational form. The solution

$$\mathbf{f}_\alpha^\dagger = \mathbf{A}_\alpha^\dagger \tilde{\mathbf{g}} = (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \tilde{\mathbf{g}} \quad (3.12)$$

is also given by the following *unconstrained* optimisation problem

$$\mathbf{f}_\alpha^\dagger = \arg \min_{\mathbf{f} \in \mathbb{R}^n} [\Phi = \|\tilde{\mathbf{g}} - \mathbf{A}\mathbf{f}\|^2 + \alpha \|\mathbf{f}\|^2] \quad (3.13)$$

More generally we could consider an optimisation problem

$$\mathbf{f}_\alpha^\dagger = \arg \min_{\mathbf{f} \in \mathbb{R}^n} [\Phi = \mathcal{D}(\tilde{\mathbf{g}}, \mathbf{A}\mathbf{f}) + \alpha \Psi(\mathbf{f})] \quad (3.14)$$

where  $\mathcal{D}$  is a data fitting functions (c.f.Eq. 2.6) and  $\Psi$  is a regularisation functional.

### 3.3.1 Generalised Tikhonov Regularisation

The particular case

$$\Psi(\mathbf{f}) = \frac{1}{2} \|\mathbf{f}\|^2 \quad (3.15)$$

is known as *zero-order Tikhonov regularisation* (the factor  $\frac{1}{2}$  is for convenience and can be absorbed into  $\alpha$ . The “zero-order” refers to the highest order of derivative in the functional, i.e. zero. Tikhonov refers in general to any *quadratic* functional

$$\Psi(\mathbf{f}) = \frac{1}{2} \|\mathbf{f}\|_{\Gamma}^2 = \frac{1}{2} \mathbf{f}^T \Gamma \mathbf{f} \quad (3.16)$$

where  $\Gamma$  is a matrix. There are two ways that this is commonly used. The first is to induce *correlation* in the elements of the solution. This is clear from the Bayesian interpretation of the regularisation functional as the negative log of a prior probability density on the solution:

$$\Psi(\mathbf{f}) = -\log \pi(\mathbf{f}) \quad \Leftrightarrow \quad \pi(\mathbf{f}) = \exp - \left[ \frac{1}{2} \mathbf{f}^T \Gamma \mathbf{f} \right] \quad (3.17)$$

where the expression on the right is a multivariate Gaussian distribution with covariance  $\mathbf{C} = \Gamma^{-1}$ . Secondly, consider that the functional is specified on the *derivative* of  $\mathbf{f}$ .

$$\Psi(\mathbf{f}) = \frac{1}{2} \|\mathbf{f}'\|^2 \quad (3.18)$$

In the discrete setting, the derivative can be implemented as a finite difference matrix giving

$$\Psi(\mathbf{f}) = \frac{1}{2} \|\mathbf{D}\mathbf{f}\|^2 = \frac{1}{2} \mathbf{f}^T \mathbf{D}^T \mathbf{D} \mathbf{f} \quad (3.19)$$

From the regularisation point of view we can argue that Eq. 3.18, is preferable to Eq. 3.15 because it imposes a penalty on the oscillations in the solution directly, rather than just a penalty on the magnitude of the solution. It is also clear that Eq. 3.19 is of the same form as Eq. 3.16 with  $\Gamma = \mathbf{D}^T \mathbf{D}$ . However, it is not conceptually straightforward to associate a covariance matrix  $\mathbf{C} = (\mathbf{D}^T \mathbf{D})^{-1}$ . This is because operator  $\mathbf{D}$  is rank-deficient : it has as a null space any constant function. From the Bayesian point of view this means that the prior is such that any solutions differing by a constant are equally probable. In fact this is desirable feature. It means that we are not biasing the solution to have a proscribed mean. Remember that we are not optimising  $\Psi(\mathbf{f})$  on its own, but in conjunction with data. It is the data that will dictate the mean of the solution.

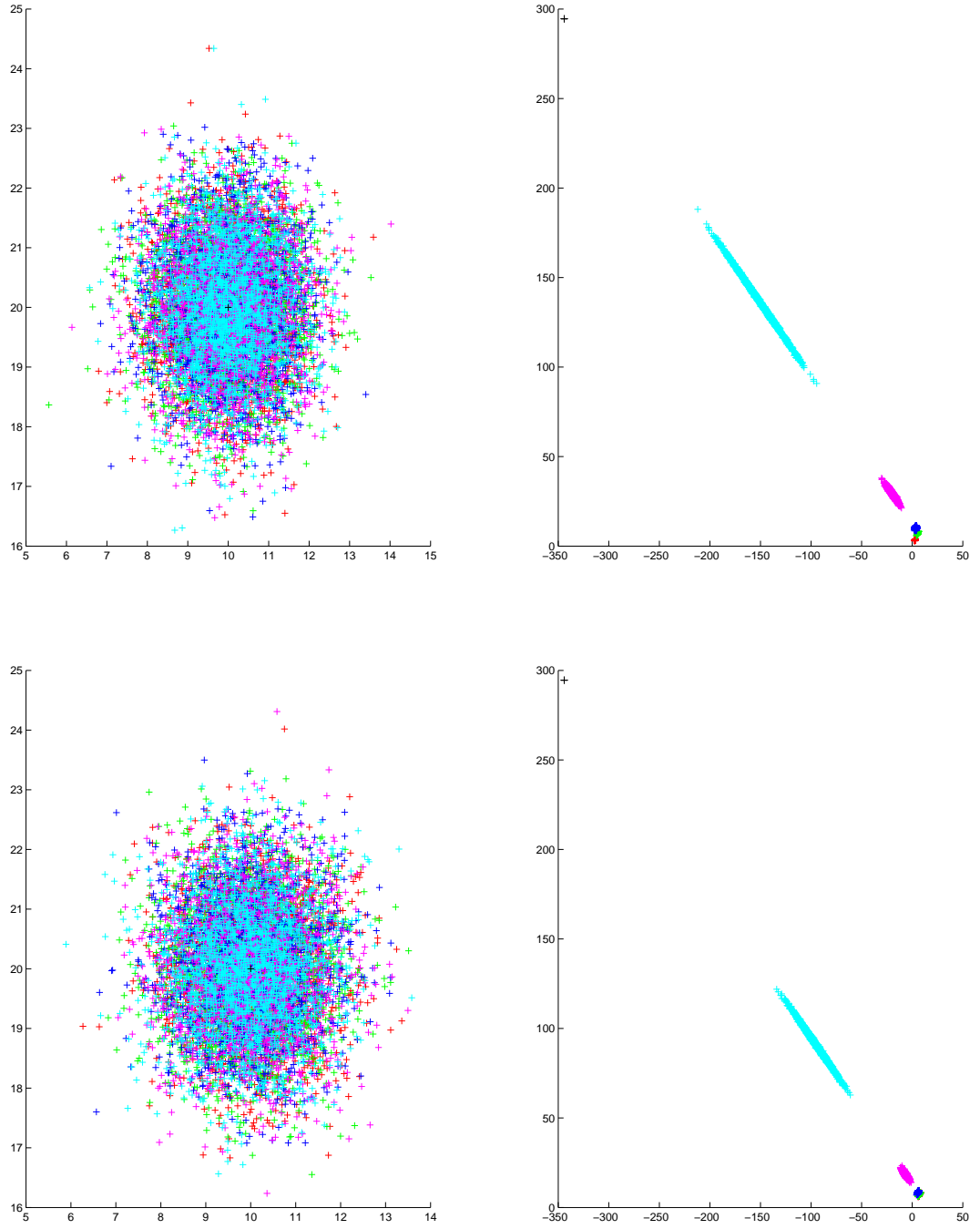


Figure 3.5: Regularised inversion of the case in figure 2.2.3; top row with prior covariance  $\mathbf{C} = \mathbf{I}$ , bottom row with  $\mathbf{C} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$ . Regularisation parameter  $\alpha = [5, 0.5, 0.05, 0.005, 0.0005]$



### 3.3.2 Total Variation

One of the disadvantages of using a Tikhonov regularisation functional is that the quadratic growth of the penalty term can be too fast, and lead to over damping of solutions that are actually desirable. For example if the true function  $\mathbf{f}$  has a step edge

$$\mathbf{f} = \begin{cases} 1 & 0.4 \leq \mathbf{x} \leq 0.6 \\ 0 & \text{otherwise} \end{cases} \quad (3.20)$$

its derivative consists of a  $\delta$ -function at  $x = 0.4$  and a negative  $\delta$ -functions at  $x = 0.6$ . Taking the square norm of a  $\delta$ -function leads to infinity in the penalty term, making such a solution impossible to achieve. In fact one can see that the form Eq. 3.18 will force the edges to be less and less steep; a desirable property to eliminate noise, but undesirable if real edges are present. If instead we use a 1-norm

$$\Psi(\mathbf{f}) = |\mathbf{f}'| \quad (3.21)$$

Then the value will be  $\Psi = 2$  for the case considered in Eq. 3.20.

The so-called *Total Variation* (TV) functional is just

$$\Psi(\mathbf{f}) = |\mathbf{D}\mathbf{f}| = \sum_{i=1}^{n-1} |f_{i+1} - f_i| \quad (3.22)$$

and will be considered in more detail for images in section 5. The use of TV as a regulariser can lead to reconstructions with sharp edges that still suppress noise. However, there is no longer a direct matrix inversion like Eq. 3.12.

## 3.4 Choice of Regularisation parameter

Choice of regularisation parameter is usually considered a tradeoff between the accuracy of recovery of the correct solution, and the accuracy of fitting to the data.

A measure of reconstruction accuracy can only be determined if a true solution is actually known. Then we may define the *estimation error*

$$\mathbf{e}_\alpha := \mathbf{f}_{\text{true}} - \mathbf{f}_\alpha^\dagger \quad (3.23)$$

Usually we only know the measured data, and the approximation to it given by the model acting on our estimated solution. We define the data residual as

$$\mathbf{r}_\alpha := \tilde{\mathbf{g}} - A\mathbf{f}_\alpha^\dagger \quad (3.24)$$

Note that the true solution  $\mathbf{f}_{\text{true}}$  should not give a zero residual, because the data contains noise. If the true solution were available we could define the *predictive error*

$$\mathbf{p}_\alpha := A\mathbf{e}_\alpha = A\mathbf{f}_{\text{true}} - A\mathbf{f}_\alpha^\dagger \quad (3.25)$$

Regularisation parameters selection methods can be divided between those where the noise statistics are known, and those where it is not. In the following we assume that the noise is identical and independent on each component of the data, i.e that it is drawn from a zero-mean Gaussian probability distribution with a multiple of identity as covariance

$$\eta \sim N(0, \sigma^2 \mathbf{I})$$

This is known as *white noise*. The case where the noise is correlated can be handled equally well by rotation of the data-space coordinates. The case where the noise is non-Gaussian is more complex and will be considered later.

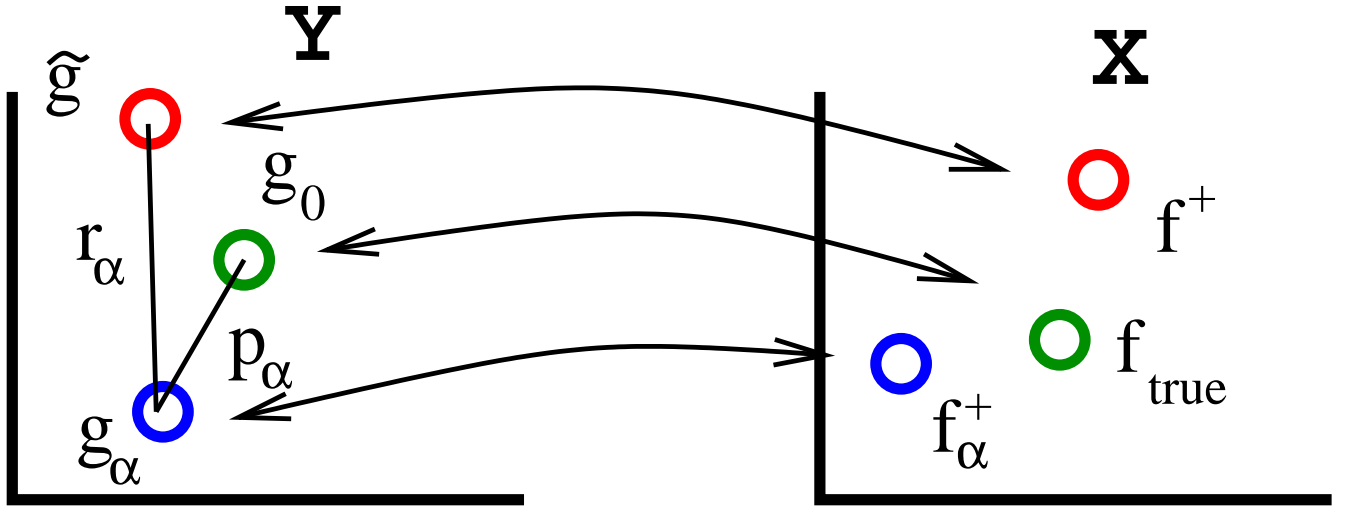


Figure 3.6: Under forward mapping the true solution  $f_{\text{true}}$  gives rise to noise free data  $g_0$ . In the presence of noise the measured data is  $\tilde{g} = g_0 + \eta$ . Reconstruction of  $\tilde{g}$  with regularisation gives a solution  $f_\alpha$  whose forward map gives modelled data  $g_\alpha$ . The *residual error* is the data-fit discrepancy  $r_\alpha$ , whereas the *predictive error*  $p_\alpha$  is the difference of the modelled data to the true, noise-free data.

### 3.4.1 Discrepancy Principle (DP)

In this method we assume that the norm of the residual should be equal to the expectation of the norm of the noise

$$\|\mathbf{r}_\alpha\|^2 = \mathbb{E}(\|\eta\|^2) = n\sigma^2 \quad (3.26)$$

where  $n$  is the length of the data vector. We therefore need to solve the nonlinear optimisation problem

$$\text{DP}(\alpha) := \frac{1}{n} \|\mathbf{r}_\alpha\|^2 - \sigma^2 = 0. \quad (3.27)$$

In some cases we can get an expression for  $DP(\alpha)$  in terms of the SVD. For example using Eq. 3.6 we can put

$$\begin{aligned} \mathbf{r}_\alpha &= (\mathbf{U}\mathbf{W}\mathbf{W}_\alpha^\dagger\mathbf{U}^\mathbf{T} - \mathbf{I}) \tilde{\mathbf{g}} \\ \Rightarrow DP(\alpha) &= \frac{1}{n} \sum_{i=1}^n (q_\alpha(w_i^2) - 1)^2 (\mathbf{u}_i \cdot \tilde{\mathbf{g}})^2 - \sigma^2 \end{aligned} \quad (3.28)$$

If  $0 \leq q_\alpha(w_i^2) < 1$  is monotonically decreasing (for example Tikhonov regularisation Eq. 3.10) then  $DP(\alpha)$  is monotonically increasing. This means that the solution to Eq. 3.27 can be found straightforwardly, e.g. by binary search.

Suppose that the noise is correlated

$$\eta \sim N(0, \mathbf{C})$$

we may write a different expression for the discrepancy principle by taking the norm of the residual with the metric  $\Gamma = \mathbf{C}^{-1}$  and similarly for the noise metric, a technique known as *pre-whitening*

$$\|\mathbf{r}_\alpha\|_\Gamma^2 = \mathbb{E}(\|\eta\|_\Gamma^2) = 1 \quad (3.29)$$

### 3.4.2 Miller Criterion

The Discrepancy Principle only takes account of the statistical properties of the *data*. Is there a similar procedure for consideration of the solution error? If we knew  $\mathbf{f}_{\text{true}}$  we could in principle require

$$\|\mathbf{e}_\alpha\|^2 - m\sigma_f^2 = 0. \quad (3.30)$$

In practice we do not know  $\mathbf{f}_{\text{true}}$  but when using a regularisation functional we are implicitly assuming that the solution is statistically distributed with probability

$$\pi(\mathbf{f}) = \exp(-\Psi(\mathbf{f})).$$

We may therefore choose a parameter  $\alpha$  so that the contributions to the total functional from the data and prior terms is about equal

$$\text{Miller}(\alpha) = \frac{\frac{1}{2}\|\mathbf{r}_\alpha\|_\Gamma^2}{\sigma^2} - \Psi(\mathbf{f}) \rightarrow 0 \quad (3.31)$$

As for the Discrepancy Principle, it is the zero of this function that is required.

### 3.4.3 Unbiased Predictive Risk Estimation (UPRE)

The UPRE is based on a statistical estimator of the mean squared norm of the predictive error

$$\frac{1}{n} \|\mathbf{p}_\alpha\|^2 = \frac{1}{n} \|A\mathbf{f}_\alpha^\dagger - A\mathbf{f}_{\text{true}}\|^2, \quad (3.32)$$

which is called the *predictive risk*. We define the  $n \times n$  *influence matrix* as

$$\mathbf{P}_\alpha = \mathbf{A}\mathbf{A}_\alpha^\dagger \quad (3.33)$$

[ Note that in terms of the SVD  $\mathbf{P}_\alpha = \mathbf{U}\mathbf{W}\mathbf{W}_\alpha^\dagger\mathbf{U}^\mathbf{T}$  ]. Using Eq. 3.3 we write

$$\mathbf{p}_\alpha = (\mathbf{P}_\alpha - \mathbf{I})\mathbf{A}\mathbf{f}_{\text{true}} + \mathbf{P}_\alpha\boldsymbol{\eta} \quad (3.34)$$

Since the first term on the right is deterministic it has the same value under expectation and the second term has expectation

$$\mathbb{E} [||\mathbf{P}_\alpha\boldsymbol{\eta}||^2] = \mathbb{E} [\boldsymbol{\eta}^\mathbf{T}\mathbf{P}_\alpha^\mathbf{T}\mathbf{P}_\alpha\boldsymbol{\eta}] = \sum_{i,j} \mathbf{P}_\alpha^\mathbf{T}\mathbf{P}_\alpha \mathbb{E} [\eta_i\eta_j] = \text{Tr}[\mathbf{P}_\alpha^\mathbf{T}\mathbf{P}_\alpha]\sigma^2 \quad (3.35)$$

We write

$$\mathbb{E} \left[ \frac{1}{n} ||\mathbf{p}_\alpha||^2 \right] = \frac{1}{n} ||(\mathbf{P}_\alpha - \mathbf{I})\mathbf{A}\mathbf{f}_{\text{true}}||^2 + \frac{\sigma^2}{n} \text{Tr}[\mathbf{P}_\alpha^\mathbf{T}\mathbf{P}_\alpha] \quad (3.36)$$

Comparing this with the Discrepancy Principle, we can write

$$\mathbf{r}_\alpha = \tilde{g} - \mathbf{A}\mathbf{f}_\alpha^\dagger = (\mathbf{P}_\alpha - \mathbf{I})\mathbf{A}\mathbf{f}_{\text{true}} + (\mathbf{P}_\alpha - \mathbf{I})\boldsymbol{\eta}, \quad (3.37)$$

so that

$$E \left[ \frac{1}{n} ||\mathbf{r}_\alpha||^2 \right] = \frac{1}{n} ||(\mathbf{P}_\alpha - \mathbf{I})\mathbf{A}\mathbf{f}_{\text{true}}||^2 + \frac{\sigma^2}{n} \text{Tr}[\mathbf{P}_\alpha^\mathbf{T}\mathbf{P}_\alpha] - \frac{2\sigma^2}{n} \text{Tr}[\mathbf{P}_\alpha] + \sigma^2, \quad (3.38)$$

and we write

$$\mathbb{E} \left[ \frac{1}{n} ||\mathbf{p}_\alpha||^2 \right] = \mathbb{E} \left[ \frac{1}{n} ||\mathbf{r}_\alpha||^2 \right] + \frac{2\sigma^2}{n} \text{Tr}[\mathbf{P}_\alpha] - \sigma^2. \quad (3.39)$$

The UPRE is the random function whos expectation is given by Eq. 3.39 :

$$\text{UPRE}(\alpha) = \frac{1}{n} ||\mathbf{r}_\alpha||^2 + \frac{2\sigma^2}{n} \text{Tr}[\mathbf{P}_\alpha] - \sigma^2. \quad (3.40)$$

The quantity required is the value of  $\alpha$  that minimises Eq. 3.40.

#### 3.4.4 Generalised Cross Validation (GCV)

The method of Generalised Cross Validation (GCV) is an alternative to UPRE that does not require prior knowledge of the variance  $\sigma^2$  of the noise. It is given by minimising

$$\text{GCV}(\alpha) = \frac{\frac{1}{n} ||\mathbf{r}_\alpha||^2}{\left( \frac{1}{n} \text{Tr}[\mathbf{I} - \mathbf{P}_\alpha] \right)^2} \quad (3.41)$$

### 3.4.5 L-curve

The L-curve is a graphical interpretation of the Miller method. To implement it, we plot a graph of the points  $\{\|\mathbf{r}_\alpha\|^2, \Psi(\mathbf{f}_\alpha^\dagger)\}$  for different values of  $\alpha$  on a log-log plot. The graph usually has a sharply defined “corner” like that of the letter “L” (hence the name). It is the value of  $\alpha$  corresponding to this corner that is chosen as the regularisation parameter value. Sometimes we fit a spline curve through a set of  $\{\|\mathbf{r}_\alpha\|^2, \Psi(\mathbf{f}_\alpha^\dagger)\}$  points, and derive the required  $\alpha$  as the point of maximum curvature of this fitted curve.

**References for this section :** [8, 5, 9].

## 4 Iterative Methods

When the dimension of the problem becomes large the SVD is an impractical tool, and direct inversion of  $A$  is often impossible. Even storage of a matrix representation of  $A$  becomes difficult if the dimensions is very large. This leads to the requirement for iterative methods. We first consider linear problems.

### 4.1 Gradient Descent and Landweber Method

Consider the minimisation of a quadratic form

$$\frac{1}{2} \|A\mathbf{f} - \mathbf{g}\|^2 \rightarrow \min \quad \equiv \quad \Phi(\mathbf{f}, \mathbf{g}) = \frac{1}{2} \langle A^T A \mathbf{f}, \mathbf{f} \rangle - \langle A^T \mathbf{g}, \mathbf{f} \rangle \rightarrow \min \quad (4.1)$$

The gradient is

$$\nabla_{\mathbf{f}} \Phi(\mathbf{f}, \mathbf{g}) = A^T A \mathbf{f} - A^T \mathbf{g} \quad (4.2)$$

which coincides with the negative residual associated with the least-squares criterion

$$\mathbf{r} = A^T(\mathbf{g} - A\mathbf{f}) \quad (4.3)$$

Given an approximation  $\mathbf{f}_k$  for the solutions, in a neighbourhood of the point the functional  $\Phi$  decreases most rapidly in the direction of the *negative* gradient, i.e. in the direction of the residual  $\mathbf{r}$ . Thus for small values  $\tau$ , we can guarantee that a new solution

$$\mathbf{f}_{k+1} = \mathbf{f}_k + \tau \mathbf{r}_k \quad (4.4)$$

will give a smaller value of  $\Phi$ . How do we choose  $\tau$ ? It makes sense to minimise the one dimensional function

$$\phi(\tau) = \Phi(\mathbf{f}_k + \tau \mathbf{r}_k, \mathbf{g}) = \Phi(\mathbf{f}_k, \mathbf{g}) + \frac{\tau^2}{2} \|A\mathbf{r}_k\|^2 - \tau \|\mathbf{r}_k\|^2. \quad (4.5)$$

Since this is a quadratic function of  $\tau$  we have immediately that

$$\tau_k = \frac{\|\mathbf{r}_k\|^2}{\|A\mathbf{r}_k\|^2}, \quad (4.6)$$

which gives rise to the *steepest descent* update rule.<sup>1</sup> We can see immediately that the residuals have the properties

- They can be obtained iteratively

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \tau_k A^T A \mathbf{r}_k, \quad (4.7)$$

---

<sup>1</sup>More generally the minimisation of the one-dimensional functional  $\phi(\tau)$  would be non-quadratic and would require a *line-search* method.

- They satisfy an orthogonality condition

$$\langle \mathbf{r}_{k+1}, \mathbf{r}_k \rangle = 0. \quad (4.8)$$

Since  $\mathbf{A}^T \mathbf{A}$  is positive-definite the method has a simple geometric interpretation. Starting at  $f_0$

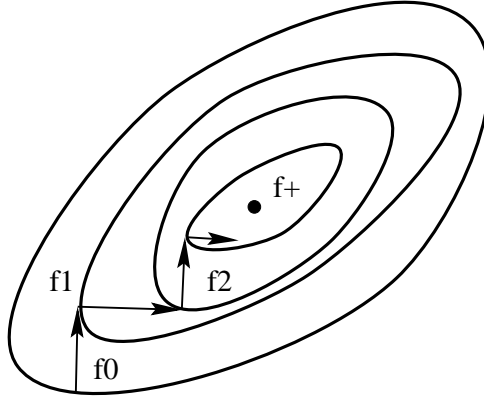


Figure 4.1: The steepest descent method.

one moves perpendicular to the level set of  $\Phi(f_0, g)$  for a distance so that at the next point  $f_1$  the line  $f_0 - f_1$  is tangent to  $\Phi(f_1, g)$ . Then one moves perpendicular to the level set of  $\Phi(f_1, g)$  for a distance so that at the next point  $f_2$  the line  $f_1 - f_2$  is tangent to  $\Phi(f_2, g)$ .

The steepest-descent method is a nonlinear method for solving least-squares problems. It works equally well for any optimisation problem sufficiently close to a local minimum. Stopping the iterations before minimisation has a regularising effect. However it can be unreasonably slow. This is clear from figure 4.1, where in the case of an anisotropic optimisation surface, the iterations can only descend very gradually to the minimum. Therefore it is almost always better to use a conjugate gradient method.

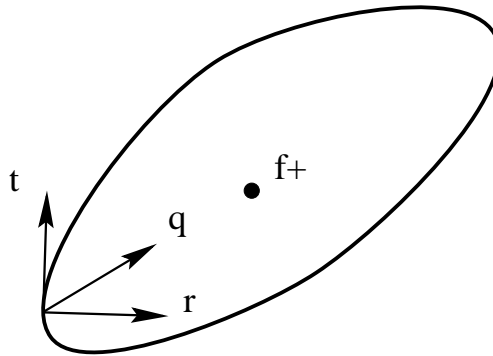


Figure 4.2: For a level surface  $\Phi = \text{constant}$ ,  $\mathbf{t}$  is a tangent vector,  $\mathbf{r}$  is the normal vector, and  $\mathbf{q}$  is the conjugate direction.

## 4.2 Conjugate Gradients

Two vectors  $\mathbf{p}$  and  $\mathbf{q}$  are said to be conjugate with respect to a matrix  $\mathbf{B}$  if

$$\langle \mathbf{p}, \mathbf{B}\mathbf{q} \rangle = 0$$

Then, given any level surface (ellipsoid) of the discrepancy functional and a point on this surface, the direction of the conjugate gradient at this point is the direction orthogonal to the tangent plane w.r.t the matrix  $\mathbf{A}^T\mathbf{A}$ .

We observe that a descent method such as described in section 4.1 builds the solution at the  $j^{\text{th}}$  iteration, from a vector space spanned by the set of vectors

$$\mathcal{K}^{(j)}(\mathbf{A}, \mathbf{g}) = \left\{ \mathbf{A}^T \mathbf{g}, \mathbf{A}^T \mathbf{A} \mathbf{A}^T \mathbf{g}, \dots (\mathbf{A}^T \mathbf{A})^{(j)} \mathbf{A}^T \mathbf{g} \right\} \quad (4.9)$$

where it is assumed that all of the vectors in the set are mutually independent. This space is called the *Krylov subspace* of  $\mathbf{A}$ . An optimal solution may be found by choosing at each iteration a direction  $\mathbf{q}$  that is the least-squares solution in this Krylov subspace. At the first iteration the subspace is of dimension one, and our only choice is the same as steepest descents. At subsequent directions we choose a direction that points towards the centre of an ellipsoid given by the intersection of the functional level surface with the Krylov subspace.

If  $\mathbf{P}_k$  is the projection operator onto  $\mathcal{K}^{(j)}(\mathbf{A}, \mathbf{g})$  then we need to solve

$$\begin{aligned} ||\mathbf{A}\mathbf{P}_k \mathbf{f} - \mathbf{g}||^2 &\rightarrow \min \\ \text{subject to} &\quad \mathbf{P}_k \mathbf{f} = \mathbf{f} \end{aligned} \quad (4.10)$$

which is given by

$$\mathbf{P}_k \mathbf{A}^T \mathbf{A} \mathbf{P}_k \mathbf{f} = \mathbf{P}_k \mathbf{A}^T \mathbf{g} \quad (4.11)$$

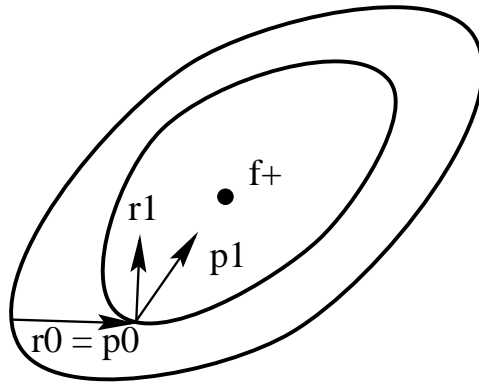


Figure 4.3: Two dimensional representation of the conjugate gradient method.

The CG method is based on the construction of two bases  $\{\mathbf{r}_k\}$  and  $\{\mathbf{p}_k\}$  as follows



- 

$$\mathbf{r}_0 = \mathbf{p}_0 = \mathbf{A}^T \mathbf{g} \quad (4.12)$$

- compute

$$\alpha_k = \frac{||\mathbf{r}_k||^2}{\langle \mathbf{r}_k, \mathbf{A}^T \mathbf{A} \mathbf{p}_k \rangle} \quad (4.13)$$

- compute

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A}^T \mathbf{A} \mathbf{p}_k \quad (4.14)$$

- compute

$$\beta_k = \frac{\langle \mathbf{r}_{k+1}, \mathbf{A}^T \mathbf{A} \mathbf{p}_k \rangle}{\langle \mathbf{p}_k, \mathbf{A}^T \mathbf{A} \mathbf{p}_k \rangle} \quad (4.15)$$

- compute

$$\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k \quad (4.16)$$

Then the iterative scheme for the solutions is

$$\begin{aligned} \mathbf{f}_0 &= 0 \\ \mathbf{f}_{k+1} &= \mathbf{f}_k + \alpha_k \mathbf{p}_k \end{aligned} \quad (4.17)$$

The coefficients  $\alpha_k, \beta_k$  are chosen so that the following orthogonality conditions hold

$$\langle \mathbf{r}_{k+1}, \mathbf{r}_k \rangle = 0, \quad \langle \mathbf{p}_{k+1}, \mathbf{A}^T \mathbf{A} \mathbf{p}_k \rangle = 0,$$

i.e. the sequence of  $\{\mathbf{r}_k\}$  are such that  $\mathbf{r}_1 \perp \mathbf{r}_0, \mathbf{r}_2 \perp \mathbf{r}_1$ , etc. and  $\mathbf{p}_1$  is conjugate to  $\mathbf{p}_0, \mathbf{p}_2$  is conjugate to  $\mathbf{p}_1$  etc. Furthermore we can show that

- the set  $\{\mathbf{r}_k\}$  form an orthogonal basis of  $\mathcal{K}^{(j)}(\mathbf{A}, \mathbf{g})$ ,
- the set  $\{\mathbf{p}_k\}$  form an  $\mathbf{A}^T \mathbf{A}$ -orthogonal basis of  $\mathcal{K}^{(j)}(\mathbf{A}, \mathbf{g})$ ,
- the set  $\{\mathbf{r}_k\}$  are the residuals

$$\mathbf{r}_k = \mathbf{A}^T \mathbf{g} - \mathbf{A}^T \mathbf{A} \mathbf{f}_k$$

- the following alternative expression for  $\alpha_k, \beta_k$  hold :

$$\alpha_k = \frac{||\mathbf{r}_k||^2}{\langle \mathbf{p}_k, \mathbf{A}^T \mathbf{A} \mathbf{p}_k \rangle} \quad \beta_k = \frac{||\mathbf{r}_{k+1}||^2}{||\mathbf{r}_k||^2}$$

which implies that  $\alpha_k, \beta_k$  are positive.

Starting at  $f_0$  one moves perpendicular to the level set of  $\Phi(f_0, g)$  for a distance so that at the next point  $f_1$  the line  $f_0 - f_1$  is tangent to  $\Phi(f_1, g)$ . Then one moves to the centre of the ellipse given by the intersection of the level surface  $\Phi(f_1, g) = \text{constant}$  with the plane spanned by  $\mathbf{r}_0$  and  $\mathbf{r}_1$ . The direction of movement is  $\mathbf{p}_1$ . Then we move to centre of the ellipse given by the

intersection of the level surface  $\Phi(f_2, g) = \text{constant}$  with the three-dimensional space spanned by  $\mathbf{r}_0, \mathbf{r}_1, \mathbf{r}_2$ . The direction of movement is  $\mathbf{p}_2$ . And so on. If the problem is  $n$ -dimensional we reach the minimum in  $n$  steps.

**References for this section :** [8, 5, 7].

## 5 Image Regularisation functionals

We here introduce a functional used frequently in image reconstruction and image processing.

$$\Psi(f) = \int_{\Omega} \psi(|\nabla f|) \, d\mathbf{x} \quad (5.1)$$

We want to derive the Gateaux derivative for Eq. 5.1 in direction  $h$ , defined to be

$$\Psi'(f)h := \lim_{\epsilon \rightarrow 0} \left[ \frac{\Psi(f + \epsilon h) - \Psi(f)}{\epsilon} \right] \quad (5.2)$$

for arbitrary function  $h$  satisfying the same boundary conditions as  $f$ .  $\Psi'(f)$  is a linear operator mapping from the space of functions to a scalar value. If we can show that it does not depend on  $h$  then  $\Psi'(f)$  is the Fréchet derivative of  $\Psi$  and constitutes a generalised notion of a gradient. See appendix section B.1 for more details.

Expanding Eq. 5.1 we have

$$\Psi(f + \epsilon h) = \int_{\Omega} \psi \left( \sqrt{\nabla f \cdot \nabla f + 2\epsilon \nabla f \cdot \nabla h + \epsilon^2 \nabla h \cdot \nabla h} \right) \, d\mathbf{x}$$

Taking the Taylor series we have

$$\Psi(f + \epsilon h) = \Psi(f) + \int_{\Omega} 2\epsilon \psi'(|\nabla f|) \frac{1}{2} |\nabla f|^{-1} \nabla f \cdot \nabla h \, d\mathbf{x} + O(\epsilon^2)$$

Making the definition

$$\kappa := \frac{\psi'(|\nabla f|)}{|\nabla f|}$$

we have

$$\Psi'(f)h = \int_{\Omega} \kappa \nabla f \cdot \nabla h \, d\mathbf{x} = - \int_{\Omega} h \nabla \cdot \kappa \nabla f \, d\mathbf{x}$$

where the last expression results from using the divergence theorem. We now define the differential operator

$$\mathcal{L}(f) = -\nabla \cdot \kappa \nabla$$

so that we can state

$$\Psi'(f)h = \int_{\Omega} h \mathcal{L}(f) f \, d\mathbf{x} = \langle h, \mathcal{L}(f) f \rangle = \langle h, \delta \Psi \rangle$$

Note that

- $\mathcal{L}(f)$  is a linear operator, but it is parameterised on  $f$  because of the function  $\kappa$ .
- $\Psi'(f)$  does not depend on  $h$ . It is a linear integral operator with kernel function  $\delta \Psi$ .
- $\delta \Psi := \mathcal{L}(f)f$  is an image.

## Examples

1. *First order Tikhonov.* If  $\psi(s) = \frac{1}{2}s^2 \Leftrightarrow \psi'(s) = s$  then  $\kappa = 1$

$$\mathcal{L}(f) = -\nabla^2$$

which is the stationary Laplacian (does not depend on  $f$ ).

2. *Total Variation.* If  $\psi(s) = s \Leftrightarrow \psi'(s) = 1$ , then  $\kappa = \frac{1}{|\nabla f|}$

$$\mathcal{L}(f) = -\nabla \cdot \frac{1}{|\nabla f|} \nabla$$

3. *Perona-Malik* If  $\psi(s) = T \log \left(1 + \frac{s}{T}\right) \Leftrightarrow \psi'(s) = \frac{sT}{T+s}$  then  $\kappa = \frac{T}{T+|\nabla f|}$

$$\mathcal{L}(f) = -\nabla \cdot \frac{1}{1 + \frac{|\nabla f|}{T}} \nabla$$

where  $T$  is a threshold

Because of the division by zero of  $\kappa = \frac{1}{|\nabla f|}$  total variation is implemented by some approximation such as  $\psi(f) = \sqrt{|\nabla f|^2 + \beta^2}$ .

Repeated application of the functional gradient generates the so-called *scale-space* of an image

$$f^{n+1} = f^n + \Psi'(f^n) \tag{5.3}$$

**References for this section :** [10, 11].

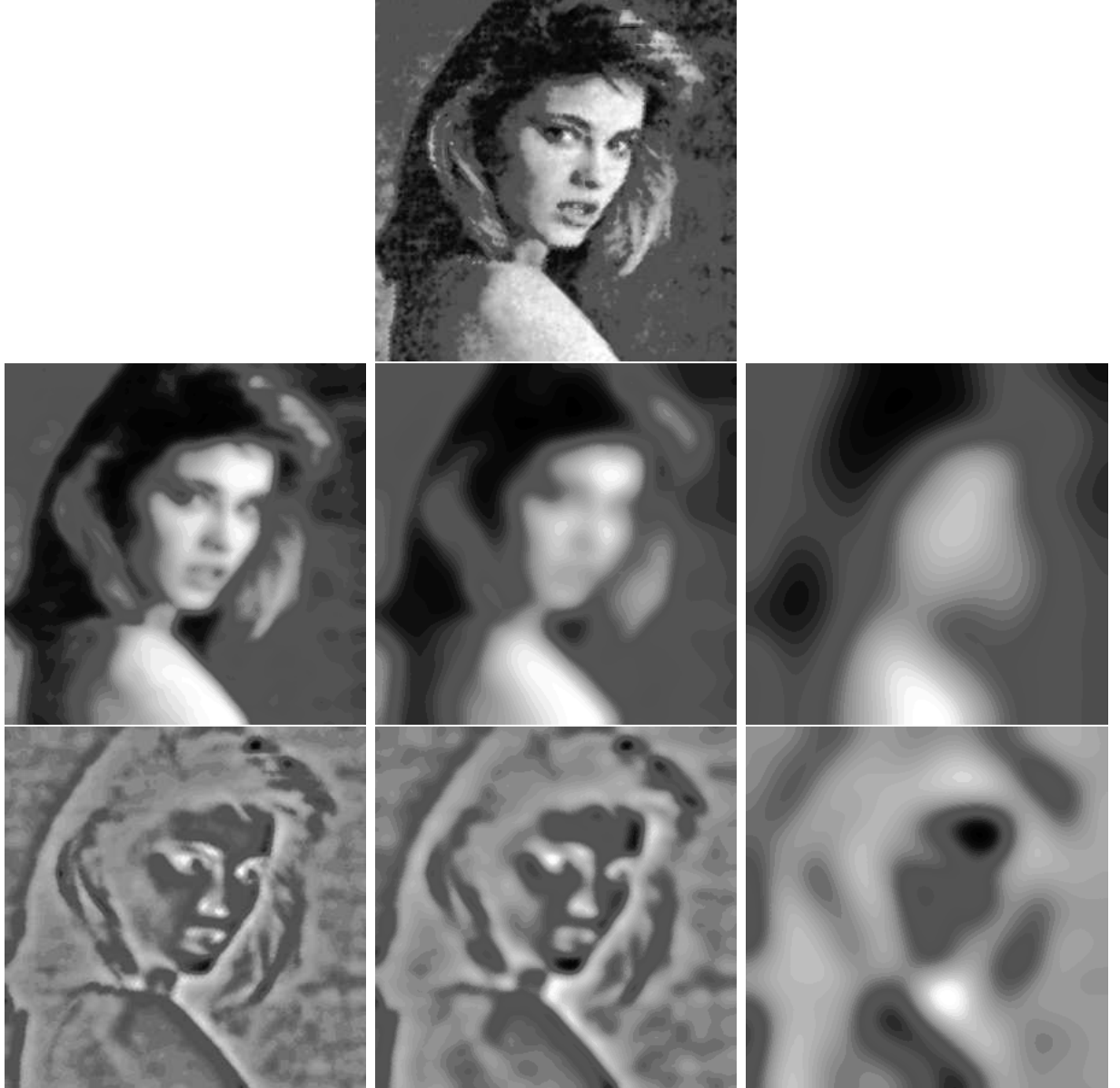


Figure 5.1: Isotropic scale resulting from the choice of a first order tikhonov prior. Iterations 1, 3, 10 corresponding to  $\Delta t = 12, 15, 40$ . Top row : original image; second row : image scale space; third row : gradient of prior  $\Psi'$ .

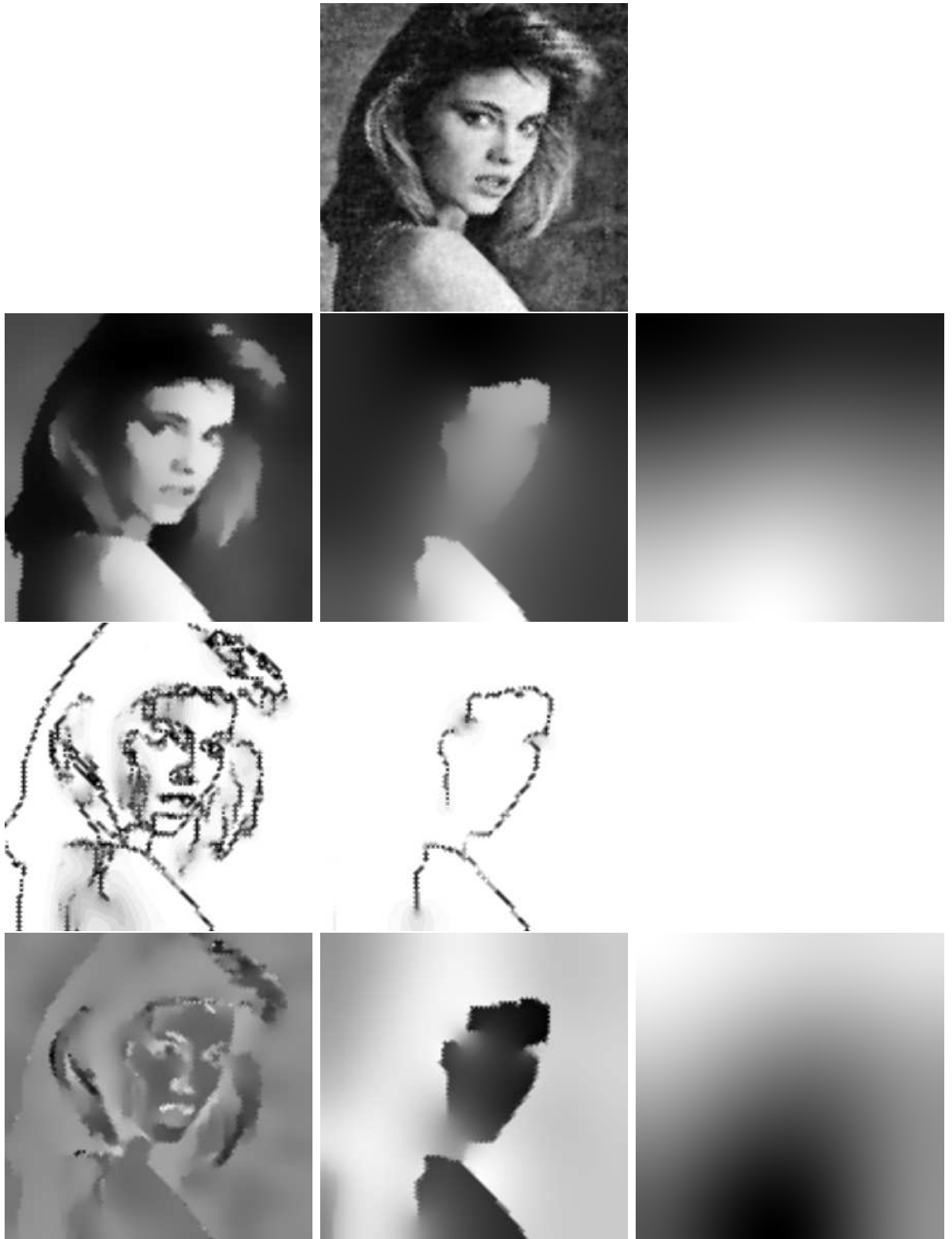


Figure 5.2: Anisotropic scale resulting from the choice of a Perona-Malik prior with threshold parameter  $T = 1.5\%$  of the image maximum. Iterations 10, 25, 35 corresponding to  $\Delta t = 40, 329, 1332$ . Top row : original image; second row : image scale space; third row : diffusivity; fourth row : gradient of prior  $\Psi'$ .

## 6 Nonlinear Optimisation

We first consider a classical example of a function that is well-posed, but strongly non-linear. The Rosenbrock function

$$\Phi(\mathbf{f}) = 100(f_2 - f_1^2)^2 + (1 - f_1)^2 \quad (6.1)$$

for which the contours are shown:

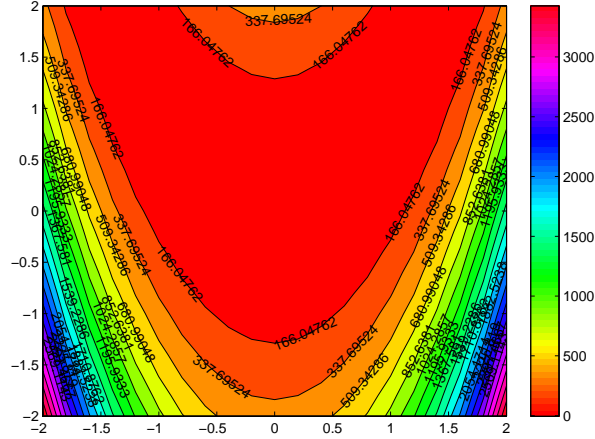


Figure 6.1: Rosenbrock function

### 6.1 Descent Methods

Both Steepest descents and Conjugate gradient have a non-linear equivalent. The key difference in each case is that i) the step length needs to be determined by a *line-search* step, and ii) a termination of *stopping* criterion needs to be determined.

We write the problem to be minimised simply as

$$\Phi(\mathbf{f}) \rightarrow \min \quad (6.2)$$

#### 6.1.1 Non-linear Steepest Descents

As an example, consider *non-linear least squares*

$$\Phi(\mathbf{f}) := \frac{1}{2} \|\mathbf{g} - A(\mathbf{f})\|^2 \quad (6.3)$$

where  $A(\mathbf{f})$  is a non-linear operator.

The algorithm is

---

**Algorithm 2** Non-linear Steepest Descent

---

$k = 0, \mathbf{f}_0 = \text{initial guess}$   
**while** stopping criterion unfulfilled **do**  
     $\mathbf{p}_k = -\Phi'(\mathbf{f}_k)$   
     $\tau_k = \arg \min_{\tau > 0} \Phi(\mathbf{f}_k + \tau \mathbf{p}_k)$   
     $\mathbf{f}_{k+1} = \mathbf{f}_k + \tau_k \mathbf{p}_k$   
     $k = k + 1$   
**end while**

---

For non-linear least squares as in Eq. 6.3 the descent direction is given by

$$-\Phi'(\mathbf{f}) = A'^*(\mathbf{f})[\mathbf{g} - A(\mathbf{f})] \quad (6.4)$$

where  $A'^*$  is the *adjoint Fréchet derivative* of operator  $A$ .

Some example stopping criteria are

- The error functional reduces by a sufficient amount in total

$$\Phi(\mathbf{f}_k) \leq \epsilon \Phi(\mathbf{f}_0)$$

- The relative change in the functional at the last iteration is smaller than some threshold

$$\frac{\Phi(\mathbf{f}_{k-1}) - \Phi(\mathbf{f}_k)}{\Phi(\mathbf{f}_{k-1})} \leq \epsilon$$

- The gradient of the functional falls below some threshold

$$\|\Phi'(\mathbf{f}_k)\| \leq \epsilon$$

The linesearch problem

$$\tau_k = \arg \min_{\tau > 0} \Phi(\mathbf{f}_k + \tau \mathbf{p}_k) \quad (6.5)$$

is a one-dimensional optimisation problem; the direction  $\mathbf{p}_k$  is called the *search direction* and the parameter  $\tau$  is called the *step length*. The search direction is a *descent direction* if

$$\langle \Phi'(\mathbf{f}_k), \mathbf{p}_k \rangle < 0. \quad (6.6)$$

Obviously, the negative gradient direction is a descent direction because

$$\langle \Phi'(\mathbf{f}_k), -\Phi'(\mathbf{f}_k) \rangle = -\|\Phi'(\mathbf{f}_k)\|^2 < 0. \quad (6.7)$$

An important result of which we will make use later is that *any symmetric positive definite operator  $B$  applied to the negative gradient direction is also a descent direction*. This follows directly from the definition of a symmetric positive definite operator :

$$\tilde{\mathbf{p}}_k = -B\Phi'(\mathbf{f}_k) \Rightarrow \langle \Phi'(\mathbf{f}_k), \tilde{\mathbf{p}}_k \rangle = \langle \Phi'(\mathbf{f}_k), -B\Phi'(\mathbf{f}_k) \rangle = -\|\Phi'(\mathbf{f}_k)\|_B^2 < 0. \quad (6.8)$$



If the problem were quadratic  $\Phi(\mathbf{f}) = \|\mathbf{A}\mathbf{f} - \mathbf{g}\|^2$ , then the linesearch would also be quadratic and would be minimised by Eq. 4.6. Instead we have to estimate the minimum. In practice it is not necessary to find the exact minimum, because the expense of doing so might not be as much as the recomputation of the gradient. If  $\mathbf{p}_k$  is a descent direction it might seem OK to take any positive step length, but there are some conditions under which either too small or too large a step could converge to a limit that was not a local minimum. Therefore the step is required to satisfy the *Wolfe conditions* :

$$\begin{aligned} \text{Wolfe condition 1 (sufficient decrease):} \quad & \phi(\tau) \leq \phi(0) + c_1\tau\phi'(0), \quad \tau > 0 \\ \text{Wolfe condition 2 (curvature condition):} \quad & \phi'(\tau) \geq c_2\tau\phi'(0), \quad \tau > 0 \\ & 0 < c_1 < c_2 < 1 \end{aligned}$$

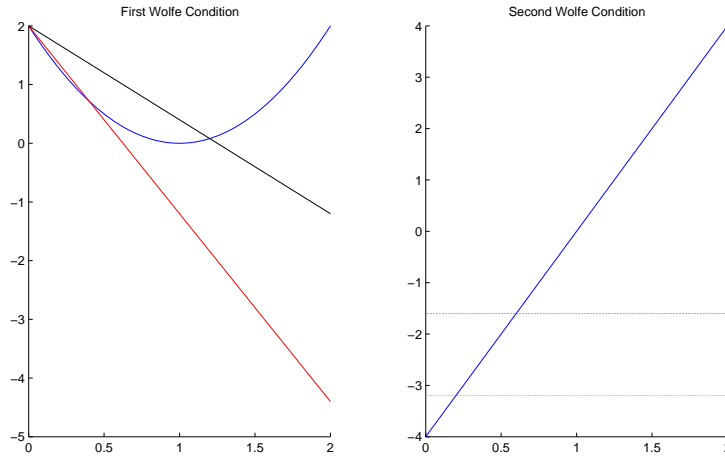


Figure 6.2: Illustration of Wolfe Conditions for simple quadratic problem.

A well-used and convenient technique for finding an approximate solution is a bracketing and quadratic interpolation technique. The aim of the bracketing part is to find three possible steps, with the middle one lower than the other two. i.e.

$$\text{find } \{\tau_a < \tau_b < \tau_c | \phi(\tau_a) > \phi(\tau_b) \& \phi(\tau_c) > \phi(\tau_b)\} \quad (6.9)$$

A plausible set can be found by a number of methods. For example if we begin with three search steps that are monotonically decreasing  $\{\tau_1 < \tau_2 < \tau_3 | \phi(\tau_1) > \phi(\tau_2) > \phi(\tau_3)\}$ , we can simply “reflect” the first point in the last :  $\tau_n = 2\tau_{n-1} - \tau_{n-3}$ , until Eq. 6.9 is satisfied.

The interpolation part simply consists of fitting a quadratic  $\phi = \alpha\tau^2 + \beta\tau + \gamma$  through three points and finding it's minimum. I.e. we need to find coefficients  $\alpha, \beta, \gamma$  so that

$$\begin{pmatrix} \phi(\tau_a) \\ \phi(\tau_b) \\ \phi(\tau_c) \end{pmatrix} = \begin{pmatrix} \tau_a^2 & \tau_a & 1 \\ \tau_b^2 & \tau_b & 1 \\ \tau_c^2 & \tau_c & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} \quad (6.10)$$

whence the required line-search minimum is given by  $\tau_{\min} = \frac{-\beta}{2\alpha}$

### 6.1.2 Non-Linear Conjugate Gradients

The non-linear CG method is derived from the CG method, again using linesearch. It can be written

---

**Algorithm 3** Non-linear Conjugate Gradients (Fletcher-Reeves)

---

```

 $k = 0, \mathbf{f}_0 = \text{initial guess}$ 
 $\mathbf{p}_0 = \mathbf{r}_0 = -\Phi'(\mathbf{f}_0)$ 
while stopping criterion unfulfilled do
     $\tau_k = \arg \min_{\tau \geq 0} \Phi(\mathbf{f}_k + \tau \mathbf{p}_k)$ 
     $\mathbf{f}_{k+1} = \mathbf{f}_k + \tau \mathbf{p}_k$ 
     $\mathbf{r}_{k+1} = -\Phi'(\mathbf{f}_{k+1})$ 
     $\beta_k = \frac{\|\mathbf{r}_{k+1}\|^2}{\|\mathbf{r}_k\|^2}$ 
     $\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$ 
     $k = k + 1$ 
end while

```

---

An alternative (more commonly used) is the Polak-Ribiere algorithm which is identical except for the step

$$\beta_k = \frac{\langle \mathbf{r}_{k+1} - \mathbf{r}_k, \mathbf{r}_{k+1} \rangle}{\langle \mathbf{r}_{k+1} - \mathbf{r}_k, \mathbf{p}_k \rangle} \quad (6.11)$$

## 6.2 Newton Methods

From a quadratic approximation of the objective function

$$\Phi(\mathbf{f} + \mathbf{h}) \simeq \Phi(\mathbf{f}) + \langle \Phi'(\mathbf{f}), \mathbf{h} \rangle + \frac{1}{2} \langle \mathbf{h}, \Phi''(\mathbf{f}) \mathbf{h} \rangle \quad (6.12)$$

if  $\Phi''(\mathbf{f})$  is positive definite then Eq. 6.12 has a unique minimiser found by solving

$$\Phi''(\mathbf{f}) \mathbf{h} = -\Phi'(\mathbf{f}), \quad (6.13)$$

and the iterative Newton scheme is

$$\mathbf{f}_{k+1} = \mathbf{f}_k - (\Phi''(\mathbf{f}_k))^{-1} \Phi'(\mathbf{f}_k) \quad (6.14)$$

Operator  $\Phi''(\mathbf{f})$  is the *Hessian* of  $\Phi$ . Note that for a linear least squares problem  $\Phi(\mathbf{f}) = \|\mathbf{A}\mathbf{f} - \mathbf{g}\|^2$  we would have :

$$-\Phi'(\mathbf{f}) \rightarrow \mathbf{A}^T(\mathbf{g} - \mathbf{A}\mathbf{f}), \quad \Phi''(\mathbf{f}) = \mathbf{A}^T \mathbf{A} \quad (6.15)$$

and Eq. 6.14 with  $\mathbf{f}_0 = \mathbf{0}$  is just the same as solving the Moore-Penrose inverse

$$\mathbf{A}^T \mathbf{A} \mathbf{f}_1 = \mathbf{A}^T \mathbf{g} \quad (6.16)$$

For a non-linear least-squares problem we have

$$-\Phi'(\mathbf{f}) \rightarrow A'^*(\mathbf{f})[\mathbf{g} - A(\mathbf{f})], \quad \Phi''(\mathbf{f}) \rightarrow A'^*(\mathbf{f})A'(\mathbf{f}) - A''(\mathbf{f})[\mathbf{g} - A(\mathbf{f})]. \quad (6.17)$$

The presence of the second term on the right in Eq. 6.17 involves the *second Fréchet derivative* of the forward mapping  $A$  which can be difficult to compute. It further can lead to the Hessian becoming no-longer positive definite. Finally, as the model fits the data better, this term tends to zero since it operates on the discrepancy  $\mathbf{g} - A(\mathbf{f})$ . For these reasons, the Newton step Eq. 6.14 is replaced by the so-called *Gauss-Newton* iteration

$$\mathbf{f}_{k+1} = \mathbf{f}_k - \left( A'^*(\mathbf{f}_k)A'(\mathbf{f}_k) \right)^{-1} A'^*(\mathbf{f}_k)[\mathbf{g} - A(\mathbf{f}_k)] \quad (6.18)$$

Note that in the Gauss-Newton, the operator  $(A'^*(\mathbf{f}_k)A'(\mathbf{f}_k))^{-1}$  is symmetric positive definite, and thus the update direction is guaranteed to be a descent direction. However, the step length given in Eq. 6.18 is not necessarily a minimiser along the direction of step, so linesearch should still be used. Thus the Newton method will be

---

**Algorithm 4** Non-linear Newton

---

```

 $k = 0, \mathbf{f}_0 = \text{initial guess}$ 
while stopping criterion unfulfilled do
   $\mathbf{p}_k = -(\Phi''(\mathbf{f}_k))^{-1} \Phi'(\mathbf{f}_k)$ 
   $\tau_k = \arg \min_{\tau > 0} \Phi(\mathbf{f}_k + \tau \mathbf{p}_k)$ 
   $\mathbf{f}_{k+1} = \mathbf{f}_k + \tau \mathbf{p}_k$ 
   $k = k + 1$ 
end while

```

---

The Gauss-Newton version is got by replacing  $\Phi''(\mathbf{f}_k)$  with its Gauss-Newton approximation  $\Phi''_{\text{GN}}(\mathbf{f}_k)$ .

### 6.2.1 Levenberg-Marquardt

An alternative to the Gauss-Newton method (which is after all, only available for problems involving an explicit mapping  $A : X \mapsto Y$ ), are methods which regularise the inversion step eq.(6.14). Specifically, if we consider

$$(\Phi''(\mathbf{f}) + \lambda I) \mathbf{h} = \Phi'(\mathbf{f}) \quad (6.19)$$

then the operator on the left will be symmetric positive definite if  $\lambda > -\min\{\sigma(\Phi''(\mathbf{f}))\}$ , i.e. if we add a diagonal term with weighting greater than the minimum of the eigenvalues of the Hessian. Then this modified Hessian will be invertible and the product with the gradient will give a descent direction. Note that this is not regularisation in the sense we have used it before. We are still only optimising the function  $\Phi(\mathbf{f})$  not an additional regularisation term, even though they have the same form.

In the *Levenberg-Marquardt* approach the weighting  $\lambda$  is variable. The direction  $\mathbf{p} = (\Phi''(\mathbf{f}) + \lambda \mathbf{I})^{-1} \Phi'(\mathbf{f})$  tends towards the descent direction, whereas as  $\lambda \rightarrow 0$  the search direction becomes closer to the Newton direction.

Starting with an initially large value, the value of  $\lambda$  is (decreased) if it proves to be a descent direction, and *increase* if not. Here is an outline of the method

---

**Algorithm 5** Levenberg-Marquardt

---

$k = 0$ ,  $\mathbf{f}_0$  = initial guess,  $\lambda$  = initial value (high)

**while** stopping criterion unfulfilled **do**

Find Step:  $\mathbf{p}_k = -(\Phi''(\mathbf{f}_k) + \lambda \mathbf{I})^{-1} \Phi'(\mathbf{f}_k)$

**if**  $\Phi(\mathbf{f}_k + \mathbf{p}_k) < \Phi(\mathbf{f}_k)$     % ( $\mathbf{p}_k$  is a descent direction)

$\lambda = \lambda / C$     % (scale down  $\lambda$ )

**else**

$\lambda = \lambda * C$     % (scale up  $\lambda$ )

**go to** Find Step

**endif**

$\mathbf{f}_{k+1} = \mathbf{f}_k + \mathbf{p}_k$     % (add the successful step. Do not use additional line search).

$k = k + 1$

**end while**

---

## 6.3 Generic methods in MatLab

Matlab offers three generic functions for optimisation :

- `fminsearch` - a method for functions when no derivatives are present
- `fminunc` - a method for unconstrained optimisation
- `fmincon` - a method for constrained optimisation

Here are some Examples

```
%-----Nedler-Mead-----
```

```
NMopts = optimset('LargeScale','off','Display','iter','TolX',1e-4, 'TolFun',1e-7,  
    'GradObj','off','PlotFcns',{@optimplotx,@optimplotfval});  
[xOptNM,funcValNM,exitFlagNM,outputStructuresNM] = fminsearch(@(x) rosenbrock(x),  
    xini,NMopts);
```

```
%-----Trust-Region-----
```

```
TRopts = optimset('LargeScale','off','Display','iter','TolX',1e-4, TolFun',1e-7,  
    'GradObj','off','PlotFcns',{@optimplotstepsize,@optimplotfval});  
[xOptTR,funcValTR,exitFlagTR,outputStructuresTR] = fminunc(@(x) rosenbrock(x),  
    xini,TRopts);
```

```
%----- Gradient Information -----
```

```
TRGopts = optimset('LargeScale','on','Display','iter','TolX',1e-4, 'TolFun',1e-7,  
    'GradObj','on','Hessian','off','PlotFcns',{@optimplotstepsize,@optimplotfval},  
[xOptTRG,funcValTRG,exitFlagTRG,outputStructuresTRG] = fminunc(@(x) rosenbrock(x),  
    xini,TRGopts);
```

```
%----- Gradient Information steepest descent -----
```

```
SDopts = optimset('LargeScale','off','Display','iter','TolX',1e-4, 'TolFun',1e-7,  
    'GradObj','on','HessUpdate','steepdesc','PlotFcns',  
    {@optimplotstepsize,@optimplotfval});  
[xOptSD,funcValSD,exitFlagSD,outputStructuresSD] = fminunc(@(x) rosenbrock(x),  
    xini,SDopts);
```

```
%----- Hessian Information -----
```

```
FNopts = optimset('LargeScale','on','Display','iter','TolX',1e-4, 'TolFun',1e-7,  
    'GradObj','on','Hessian','on','PlotFcns',  
    {@optimplotstepsize,@optimplotfval});  
[xOptFN,funcValFN,exitFlagFN,outputStructuresFN] = fminunc(@(x) rosenbrockFN(x),  
    xini,FNopts);
```



## 7 Statistical Methods

So far we considered linear and non-linear least-squares problems. The key point was that the *data-fitting* metric was a squared error. From the Bayesian point of view this corresponds to saying that the measurement error is Gaussian distributed with zero mean; the case of a non-diagonal data covariance matrix is also included in this model. However, there are other models for the measurement noise.

### 7.1 Maximum Likelihood

The Maximum likelihood method is the probabilistic interpretation of data fitting. Suppose a random vector  $\mathbf{X}$  has a joint probability density function  $P_{\mathbf{X}}(\mathbf{x}; \boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is an unknown parameter vector, and a data vector  $\mathbf{g}$  is a given realisation of  $\mathbf{X}$ , i.e. a draw from the probability space with distribution  $P_{\mathbf{X}}$ .

A *maximum likelihood estimator* (MLE) for  $\boldsymbol{\theta}$  given  $\mathbf{g}$  is the vector  $\hat{\boldsymbol{\theta}}$  which maximises the *likelihood* function

$$L(\boldsymbol{\theta}) := P_{\mathbf{X}}(\mathbf{g}; \boldsymbol{\theta}) \quad (7.1)$$

More conveniently it is the minimiser of the *negative log likelihood*

$$\ell(\boldsymbol{\theta}) := -\log P_{\mathbf{X}}(\mathbf{g}; \boldsymbol{\theta}) \quad (7.2)$$

Now we suppose that the forward model gives rise to observables

$$\mathbf{y} = A(\mathbf{f}) + \boldsymbol{\eta} \quad (7.3)$$

that are random variables by virtue of the  $\boldsymbol{\eta}$  being a random variable. We assume that we can find the parameters  $\boldsymbol{\theta}$  of the PDF of the data by taking moments of the forward model :

$$\boldsymbol{\theta} = \mathbb{E}[q(\mathbf{y})]$$

for some function  $q$ . E.g.

$$\begin{aligned} \bar{\boldsymbol{\theta}} &= \mathbb{E}[\mathbf{y}] \\ \text{cov}(\boldsymbol{\theta}) &= \mathbb{E}[(\mathbf{y} - \bar{\boldsymbol{\theta}})(\mathbf{y} - \bar{\boldsymbol{\theta}})^T] \end{aligned}$$

#### 7.1.1 Gaussian Probability Model

A continuous random vector  $\mathbf{X}$  has a Gaussian, or *Normal* distribution if its multivariate PDF is

$$P(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C}) = \frac{1}{\sqrt{(2\pi)^n \det \mathbf{C}}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (7.4)$$

The Normal distribution has  $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$  and  $\text{cov}(\mathbf{X}) = \mathbf{C}$  and is denoted

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{C})$$

If the noise is modelled by a Normal distribution, using Eq. 7.4 we can find the *maximum likelihood estimate* for  $\boldsymbol{\mu}$  given data  $\mathbf{g}$  by maximising  $P(\mathbf{g}; \boldsymbol{\mu})$ , equivalent to *minimising* the negative loglikelihood

$$\ell(\boldsymbol{\mu}) = -\log P(\mathbf{g}; \boldsymbol{\mu}, \mathbf{C}) = \frac{1}{2}(\mathbf{g} - \boldsymbol{\mu})^T \mathbf{C}^{-1}(\mathbf{g} - \boldsymbol{\mu}) + k = \frac{1}{2}\|\mathbf{g} - \boldsymbol{\mu}\|_{\mathbf{C}}^2 + k \quad (7.5)$$

where  $k$  is a constant independent of  $\boldsymbol{\mu}$ . We say that

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{A}(\mathbf{f})] = \mathbf{A}(\mathbf{f}), \quad (7.6)$$

since  $E[\boldsymbol{\eta}] = \mathbf{0}$ . We thus find the minimiser of  $\ell(\boldsymbol{\mu})$  is given by minimising

$$\Phi(\mathbf{f}, \mathbf{g}) = \frac{1}{2}(\mathbf{g} - \mathbf{A}(\mathbf{f}))^T \mathbf{C}^{-1}(\mathbf{g} - \mathbf{A}(\mathbf{f})) \quad (7.7)$$

A Newton iteration would give

$$\mathbf{f}_{k+1} = \mathbf{f}_k + \tau_k \left( \mathbf{A}'^*(\mathbf{f}_k) \mathbf{C}^{-1} \mathbf{A}'(\mathbf{f}_k) \right)^{-1} \mathbf{A}'^*(\mathbf{f}_k) \mathbf{C}^{-1} [\mathbf{g} - \mathbf{A}(\mathbf{f}_k)] \quad (7.8)$$

### 7.1.2 Poisson Probability Model

A discrete random vector  $\mathbf{X}$  has a Poisson distribution with independent components if its multivariate PDF is

$$P(\mathbf{x}; \boldsymbol{\lambda}) = \begin{cases} \prod_{i=1}^n \frac{e^{-\lambda_i} \lambda_i^{x_i}}{x_i!}, & x_i \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (7.9)$$

The Poisson distribution has  $\mathbb{E}(\mathbf{X}) = \boldsymbol{\lambda}$  and  $\text{cov}(\mathbf{X}) = \text{diag}(\boldsymbol{\lambda})$  and is denoted

$$\mathbf{X} \sim \text{Poisson}(\boldsymbol{\lambda})$$

If the noise is modelled by a Poisson distribution, using Eq. 7.9 we can find the maximum likelihood estimate for  $\boldsymbol{\lambda}$  given data  $\mathbf{g}$  by maximising  $P(\mathbf{g}; \boldsymbol{\lambda})$ , equivalent to minimising the negative loglikelihood

$$\ell(\boldsymbol{\lambda}) = \sum_{i=1}^N (\lambda_i - g_i \log \lambda_i) + k \quad (7.10)$$

where  $k$  is a constant independent of  $\boldsymbol{\lambda}$ . Now we suppose that the data is obtained through a forward operator, so that the actual Poisson variables are represented as

$$\boldsymbol{\lambda} = \mathbf{A}(\mathbf{f})$$



First take  $\mathbf{A}$  to be linear. Eq. 7.10 becomes

$$\ell(\mathbf{g}, \mathbf{A}\mathbf{f}) = \sum_{i=1}^N \left( \sum_{j=1}^M A_{ij} f_j - g_i \log \left[ \sum_{j=1}^M A_{ij} f_j \right] \right) + k \quad (7.11)$$

from which we have

$$\frac{\partial \ell(\mathbf{g}, \mathbf{A}\mathbf{f})}{\partial \mathbf{f}} = \mathbf{A}^T \mathbf{1} - \mathbf{A}^T \left[ \frac{\mathbf{g}}{\mathbf{A}\mathbf{f}} \right] \quad (7.12)$$

where  $\mathbf{1}$  is a vector of all ones, the same length as the data. To ensure positivity, make a change of variables  $\mathbf{f} = \mathbf{e}^\zeta$ , whence

$$\frac{\partial \ell}{\partial \zeta} = \frac{\partial \ell}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \zeta} = \frac{\partial \ell}{\partial \mathbf{f}} \mathbf{f}$$

which leads to

$$\mathbf{f} \odot \mathbf{A}^T \mathbf{1} = \mathbf{f} \odot \mathbf{A}^T \left[ \frac{\mathbf{g}}{\mathbf{A}\mathbf{f}} \right] \quad (7.13)$$

where  $\odot$  implies point wise multiplication. Now consider the operator

$$M(\mathbf{f}) = \frac{1}{\mathbf{A}^T \mathbf{1}} \mathbf{f} \odot \mathbf{A}^T \left[ \frac{\mathbf{g}}{\mathbf{A}\mathbf{f}} \right] \quad (7.14)$$

Then the maximum likelihood solution is a *fixed point* of this operator, i.e.

$$\mathbf{f}_{\text{ML}} = M(\mathbf{f}_{\text{ML}}) \quad (7.15)$$

The iterative *Expectation Maximisation* algorithm (called *Richardson-Lucy* in image processing) finds this fixed-point by iteration:

$$\mathbf{f}_{k+1} = \frac{\mathbf{f}_k}{\mathbf{A}^T \mathbf{1}} \odot \mathbf{A}^T \left[ \frac{\mathbf{g}}{\mathbf{A}\mathbf{f}_k} \right] \quad (7.16)$$

The above derivation is a little unrigorous. We can derive the same algorithm using rigorous EM theory, and using constrained optimisation theory, as we will see later.

In the non-linear case, the minimum of Eq. 7.10 requires

$$A'^*(\mathbf{f}) \mathbf{1} = A'^*(\mathbf{f}) \frac{\mathbf{g}}{A(\mathbf{f})} \quad (7.17)$$

and we obtain the iteration

$$\mathbf{f}_{k+1} = \frac{\mathbf{f}_k}{A'^*(\mathbf{f}_k) \mathbf{1}} \odot A'^*(\mathbf{f}_k) \frac{\mathbf{g}}{A(\mathbf{f}_k)} \quad (7.18)$$

### 7.1.3 Gaussian Approximation to a Poisson Noise Process

The Central Limit Theorem tells us that in the limit of “large enough” number of observations all probability distributions tend to the Normal distribution.

In the case of a Poisson process with rate  $\lambda$  the limiting Normal distribution is one with mean  $\lambda$  and variance  $\lambda$ . Thus a possible alternative to minimising the negative Poisson log likelihood Eq. 7.10 is to minimise

$$\Phi(\mathbf{f}, \mathbf{g}) = \frac{1}{2}(\mathbf{g} - \mathbf{A}\mathbf{f})^T (\text{diag}[\mathbf{A}\mathbf{f}])^{-1} (\mathbf{g} - \mathbf{A}\mathbf{f}) \equiv \frac{1}{2} \left\| \frac{(\mathbf{g} - \mathbf{A}\mathbf{f})}{\sqrt{\mathbf{A}\mathbf{f}}} \right\|^2 \quad (7.19)$$

This is highly non-linear. We get another approximation if we use the data as an estimator of the variance and thus use

$$\Phi(\mathbf{f}, \mathbf{g}) = \frac{1}{2}(\mathbf{g} - \mathbf{A}\mathbf{f})^T (\text{diag}[\mathbf{g}])^{-1} (\mathbf{g} - \mathbf{A}\mathbf{f}) \equiv \frac{1}{2} \left\| \frac{(\mathbf{g} - \mathbf{A}\mathbf{f})}{\sqrt{\mathbf{g}}} \right\|^2 \quad (7.20)$$

which is solved (for a linear operator ) by

$$\mathbf{f}_* = \left( \mathbf{A}^* \text{diag} \left[ \frac{1}{\mathbf{g}} \right] \mathbf{A} \right)^{-1} \mathbf{A}^* \mathbf{1} \quad (7.21)$$

We can also use a gradient method,

$$\mathbf{f}^{(n+1)} = \mathbf{f}^{(n)} + \tau \left( \mathbf{A}^* \text{diag} \left[ \frac{1}{\mathbf{g}} \right] \mathbf{A} \right)^{-1} \mathbf{A}^* \left[ \mathbf{1} - \frac{\mathbf{A}\mathbf{f}^{(n)}}{\mathbf{g}} \right] \quad (7.22)$$

## 7.2 Kullback-Liebler and other divergences

In information theory, metrics are sometimes defined without a specific interpretation in terms of measurement noise

The Kullback-Liebler information divergence is an information-theoretic measure of the relative entropy between two functions  $u, v$

$$KL(u, v) = \int \left( u(x) \log \left[ \frac{u(x)}{v(x)} \right] - u(x) + v(x) \right) dx. \quad (7.23)$$

If the functions are probability distributions the expression simplifies to the *relative entropy*

$$KL(P_1, P_2) = \int P_1(x) \log \left[ \frac{P_1(x)}{P_2(x)} \right] dx \quad (7.24)$$

since  $\int P_1 = \int P_2 = 1$ . The KL-divergence is always positive unless  $u = v$  when  $KL(u, u) = 0$ . Note that the measure is not symmetric i.e.  $KL(u, v) \neq KL(v, u)$ .

The discrete version of KL-divergence measures the difference between vectors

$$KL(\mathbf{u}, \mathbf{v}) = \sum_i \left( u_i \log \left[ \frac{u_i}{v_i} \right] - u_i + v_i \right). \quad (7.25)$$

Now consider

$$\Phi(\mathbf{f}, \mathbf{g}) = KL(\mathbf{g}, A(\mathbf{f})) \Rightarrow \Phi'(\mathbf{f}) = A'^*(\mathbf{f}) \left[ \mathbf{1} - \frac{\mathbf{g}}{A(\mathbf{f})} \right] \quad (7.26)$$

This is the same as the gradient of the negative Poisson log likelihood (NPLL) given in Eq. 7.17. Note that both the KL-divergence and its derivative are zero for  $\mathbf{g} = A(\mathbf{f})$ , whereas the NPLL is not necessarily zero (but its derivative is); i.e. the NPLL and KL divergence differ by a constant.

The Hessian of the KL-divergence is

$$\Phi''(\mathbf{f}) = A'^*(\mathbf{f}) \text{diag} \left( \frac{\mathbf{g}}{(A(\mathbf{f}))^2} \right) A'(\mathbf{f}) \quad (7.27)$$

which is symmetric positive definite.

### 7.3 The Bayesian MAP estimate

ML methods suffer from ill-posedness in just the same way that least-squares problems do. Rather than maximising the likelihood we maximise the Bayesian posterior

$$P(\boldsymbol{\theta} | \mathbf{X}) = \frac{P(\mathbf{X} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})}{P(\mathbf{X})} \quad (7.28)$$

here  $P(\mathbf{X} | \boldsymbol{\theta}) \equiv L(\mathbf{X}; \boldsymbol{\theta})$ . The maximiser of Eq. 7.28 is the minimiser of

$$\Phi(\mathbf{f}, \mathbf{g}) = \ell(\mathbf{g}; A(\mathbf{f})) + \Psi(\mathbf{f}) \quad (7.29)$$

where  $\Psi(\mathbf{f}) = -\log \pi(\mathbf{f})$ . We find the gradient

$$\Phi'(\mathbf{f}, \mathbf{g}) = A'^*(\mathbf{f}) \left[ \mathbf{1} - \frac{\mathbf{g}}{A(\mathbf{f})} \right] + \Psi'(\mathbf{f}) \quad (7.30)$$

and we can apply the same derivation as in section 7.1.2 to obtain the *One Step Late* EM algorithm as

$$\mathbf{f}_{k+1} = \frac{\mathbf{f}_k}{A'^*(\mathbf{f}_k) \mathbf{1} + \Psi'(\mathbf{f})} \odot A'^*(\mathbf{f}_k) \frac{\mathbf{g}}{A(\mathbf{f}_k)} \quad (7.31)$$

This algorithm adds the regularisation term in the normalisation term, and now, obviously, it is impossible for the solution to exactly solve the original problem  $\mathbf{g} = A(\mathbf{f})$  because this does not cancel the ratio

$$\frac{A'^*(\mathbf{f}_k) \mathbf{1}}{A'^*(\mathbf{f}_k) \mathbf{1} + \Psi'(\mathbf{f})};$$

however, this is not a contradiction : in the presence of regularisation we do not *want* the model to fit the data exactly, because doing so would propagate noise into the solution. This is the whole rationale for regularisation.

There is a danger with Eq. 7.31 : The gradient of the regulariser is often a derivative operator and may cause the denominator to become negative. Thus positivity might be lost.

## 7.4 Entropy as a prior

We can use Information measures as a prior. For example, the entropic prior is

$$\Psi(\mathbf{f}) = Ent(\mathbf{f}) = \sum_{i=1}^n f_i \log f_i - f_i \quad (7.32)$$

with derivative

$$\Psi'(\mathbf{f}) = Ent'(\mathbf{f}) = \log \mathbf{f} \quad (7.33)$$

Alternative, the KL-divergence can be used to provide the cross-entropy of *Mutual Information* to a reference vector  $\mathbf{v}$

$$\Psi(\mathbf{f}, \mathbf{v}) = KL(\mathbf{f}, \mathbf{v}) = \sum_{i=1}^n (f_i \log f_i - f_i) - (f_i \log v_i - v_i) \quad (7.34)$$

with derivative

$$\Psi'(\mathbf{f}) = KL'(\mathbf{f}) = \log \left[ \frac{\mathbf{f}}{\mathbf{v}} \right] \quad (7.35)$$

## 8 Constrained Optimisation

A general constrained optimisation problem is stated

$$\text{minimise} \quad \Phi(\mathbf{f}) \quad (8.1)$$

$$\text{subject to} \quad c_i(\mathbf{f}) \leq 0, \quad i = 1 \dots L \quad (8.2)$$

$$\text{and} \quad c_e(\mathbf{f}) = 0, \quad e = 1 \dots M \quad (8.3)$$

where  $c_i$  represent  $L$  inequality constraints and  $c_e$  represent  $M$  equality constraints.

### 8.1 Parameter Transformation

Before looking at a general solution, we can consider a commonly occuring constraint and a simple work round. If  $\mathbf{f} \succeq 0$  (i.e.  $f_i \geq 0 \forall i$ ) we say that  $\mathbf{f}$  requires a *positivity constraint*. We can make a simple change of variables

$$\boldsymbol{\zeta} = \log \mathbf{f} \quad \Leftrightarrow \quad \mathbf{f} = e^{\boldsymbol{\zeta}} \quad (8.4)$$

and then solve an unconstrained problem for  $\boldsymbol{\zeta}$  followed by exponentiation. For a non-linear problem we would have

$$\begin{aligned} A(\mathbf{f}) &\rightarrow A(\boldsymbol{\zeta}) = A(e^{\boldsymbol{\zeta}}) \\ \Rightarrow \quad A'(\boldsymbol{\zeta}) &= A'(\mathbf{f}) \circ \frac{\partial \mathbf{f}}{\partial \boldsymbol{\zeta}} = A'(\mathbf{f}) \circ \mathbf{f} \end{aligned}$$

Finding an update  $\mathbf{h}_{\boldsymbol{\zeta}}$  we could then write the update in  $\mathbf{f}$  as a multiplicative update

$$\boldsymbol{\zeta}_{n+1} = \boldsymbol{\zeta}_n + \mathbf{h}_{\boldsymbol{\zeta}} \quad \Leftrightarrow \quad \mathbf{f}_{n+1} = \mathbf{f}_n \odot e^{\mathbf{h}_{\boldsymbol{\zeta}}}$$

We could have used other transformations in place of Eq. 8.4. Another useful choice is the *logit Transform*

$$\boldsymbol{\zeta} = \log \frac{\mathbf{f} - a}{\mathbf{f} - b} \quad \Leftrightarrow \quad \mathbf{f} = \frac{be^{\boldsymbol{\zeta}} - a}{e^{\boldsymbol{\zeta}} - 1} \quad (8.5)$$

which constrains  $\mathbf{f}$  to lie between bounds  $a \preceq \mathbf{f} \preceq b$ .

### 8.2 Lagrange Multipliers and Dual Variables

Lagrangian is

$$\mathcal{L}(\mathbf{f}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \Phi(\mathbf{f}) + \sum_{i=1}^L \lambda_i c_i(\mathbf{f}) + \sum_{e=1}^M \mu_e c_e(\mathbf{f}) \quad (8.6)$$

The *primary variables* are  $\mathbf{f}$  and the *dual variables* are  $\boldsymbol{\lambda}, \boldsymbol{\mu}$

### 8.3 Karush-Kuhn-Tucker (KKT) conditions

KKT conditions are a generalisation of Lagrange variables to inequalities.

The KKT conditions have to be satisfied at the optimal solution  $\{\mathbf{f}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$

$$c_i(\mathbf{f}^*) \leq 0, \quad i = 1 \dots L \quad (8.7)$$

$$c_e(\mathbf{f}^*) = 0, \quad e = 1 \dots M \quad (8.8)$$

$$\lambda_i^* \geq 0, \quad i = 1 \dots L \quad (8.9)$$

$$\lambda_i^* c_i(\mathbf{f}^*) = 0, \quad i = 1 \dots L \quad (8.10)$$

$$\frac{\partial \Phi(\mathbf{f})}{\partial \mathbf{f}}|_{\mathbf{f}^*} + \sum_{i=1}^L \lambda_i \frac{\partial c_i(\mathbf{f})}{\partial \mathbf{f}}|_{\mathbf{f}^*} + \sum_{e=1}^M \mu_e \frac{\partial c_e(\mathbf{f})}{\partial \mathbf{f}}|_{\mathbf{f}^*} = 0 \quad (8.11)$$

### 8.4 Quadratic Programming

#### 8.4.1 Equality Constraints

Consider the constrained optimisation problem

$$\text{minimise} \quad \frac{1}{2} \mathbf{f}^T \mathbf{P} \mathbf{f} + \mathbf{q}^T \mathbf{f} + r \quad (8.12)$$

$$\text{subject to} \quad \mathbf{A} \mathbf{f} = \mathbf{b} \quad (8.13)$$

where  $\mathbf{P}$  is a symmetric positive definite  $N \times N$  matrix,  $\mathbf{q}, \mathbf{f}$  are vectors of length  $N$ ,  $\mathbf{A}$  is a  $M \times N$  constraint matrix (with  $M < N$ ) and  $\mathbf{b}$  is a vector of length  $M$ .

It is equality only problem. The Lagrangian is written

$$\mathcal{L}(\mathbf{f}, \boldsymbol{\mu}) = \frac{1}{2} \mathbf{f}^T \mathbf{P} \mathbf{f} + \mathbf{q}^T \mathbf{f} + r + \boldsymbol{\mu}^T [\mathbf{A} \mathbf{f} - \mathbf{b}] \quad (8.14)$$

We see that

$$\frac{\partial \mathcal{L}(\mathbf{f}, \boldsymbol{\mu})}{\partial \mathbf{f}} = \mathbf{P} \mathbf{f}^* + \mathbf{q} + \mathbf{A}^T \boldsymbol{\mu}^* \quad (8.15)$$

$$\frac{\partial \mathcal{L}(\mathbf{f}, \boldsymbol{\mu})}{\partial \boldsymbol{\mu}} = \mathbf{A} \mathbf{f}^* - \mathbf{b} \quad (8.16)$$

Only the 2nd and 5th KKT conditions are relevant. These are exactly what we see in eqs.(8.15) and eqs.(8.16)

$$\mathbf{A} \mathbf{f}^* = \mathbf{b}, \quad \mathbf{P} \mathbf{f}^* + \mathbf{q} + \mathbf{A}^T \boldsymbol{\mu}^* = 0$$

which can be written as

$$\begin{pmatrix} \mathbf{P} & \mathbf{A}^T \\ \mathbf{A} & 0 \end{pmatrix} \begin{pmatrix} \mathbf{f}^* \\ \boldsymbol{\mu}^* \end{pmatrix} = \begin{pmatrix} -\mathbf{q} \\ \mathbf{b} \end{pmatrix}$$

Since  $\mathbf{P}$  was spd, so is the augmented matrix

### 8.4.2 Inequality Constraints and the Active Set Method

Now consider the more general quadratic programming problem, with linear constraints

$$\text{minimise} \quad \frac{1}{2} \mathbf{f}^T \mathbf{P} \mathbf{f} + \mathbf{q}^T \mathbf{f} + r \quad (8.17)$$

$$\text{subject to} \quad \mathbf{a}_i^T \mathbf{f} \geq b_i, \quad i = 1 \dots L \quad (8.18)$$

$$\text{and} \quad \mathbf{a}_e^T \mathbf{f} = b_e, \quad e = 1 \dots M \quad (8.19)$$

In the *primal active set method* certain constraints, indexed by the *active set*  $\mathcal{A}$ , are regarded as equalities whilst the rest are temporarily disregarded. On each iteration the constraints that are active are reevaluated. If  $\mathbf{f}$  is feasible and  $\mathbf{a}_i^T \mathbf{f} = b_i$  we say that the  $i^{\text{th}}$  inequality constraint is *active* at  $\mathbf{f}$ . If  $\mathbf{a}_i^T \mathbf{f} > b_i$  then this constraint is *inactive*. The equality constraints are always active. Note the connection to the KKT condition eq(8.10) : when the  $i^{\text{th}}$  inequality constraint is satisfied, the corresponding Lagrangian variable  $\lambda_i$  must be zero and it is removed from the active set. When the inequality is not satisfied the  $\lambda_i$  has a positive value and the system seeks to find a point on the feasible region boundary for this constraint.

Now consider that  $\mathbf{f}$  is a feasible point (i.e. satisfying the constraints) and notice that

$$\Phi(\mathbf{f} + \mathbf{h}) = \Phi(\mathbf{f}) + \mathbf{h}^T \mathbf{q} + \frac{1}{2} \mathbf{h}^T \mathbf{P} \mathbf{h}$$

Also since for each active constraint  $\mathbf{a}_i^T \mathbf{f} = b_i$  then

$$\mathbf{a}_i^T (\mathbf{f} + \mathbf{h}) = b_i \Rightarrow \mathbf{a}_i^T \mathbf{h} = 0$$

This means that an update step  $\mathbf{h}$  must be tangent to the hyperplane of each constraint  $i$ . Therefore, given a feasible point  $\mathbf{f}$  and a set of active constraints, we can solve the equality constrained quadratic equation

$$\text{minimise} \quad \frac{1}{2} \mathbf{h}^T \mathbf{P} \mathbf{h} + \mathbf{q}^T \mathbf{h} \quad (8.20)$$

$$\text{subject to} \quad \mathbf{a}_i^T \mathbf{h} = 0, \quad i \in \mathcal{A} \quad (8.21)$$

$$\text{and} \quad \mathbf{a}_e^T \mathbf{h} = 0, \quad e = 1 \dots M \quad (8.22)$$

which results to solve a linear problem

$$\begin{pmatrix} \mathbf{P} & \tilde{\mathbf{A}}_I^T & \mathbf{A}_E^T \\ \tilde{\mathbf{A}}_I & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_E & \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{h} \\ \boldsymbol{\lambda} \\ \boldsymbol{\mu} \end{pmatrix} = \begin{pmatrix} -\mathbf{q} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \quad (8.23)$$

where  $\tilde{\mathbf{A}}_I$  is the matrix of the inequality constraints that are active. The vector  $\mathbf{h}$  is an update for  $\mathbf{f}$  that satisfies the currently active constraints. It needs to be tested also against the other, currently inactive constraints. If it satisfies these then

$$\mathbf{f} \rightarrow \mathbf{f} + \mathbf{h}$$

otherwise, a linesearch is taken in direction  $\mathbf{h}$

$$\tau = \min \left[ 1, \min_{j \in \mathcal{A} \& (\mathbf{a}_j^T \mathbf{h} > 0)} \frac{b_j - \mathbf{a}_j^T \mathbf{f}}{\mathbf{a}_j^T \mathbf{h}} \right] \quad (8.24)$$

and

$$\mathbf{f} \rightarrow \mathbf{f} + \tau \mathbf{h}$$

Note that eq(8.24) gives the projection onto the hyperplanes of all the inactive constraints. The minimum of the projection lengths is the closet hyperplane. Since only the inactive constraints are projected, the projection distances are always positive. If  $\tau < 1$  then a new index is added to the active set list corresponding to the index  $p$  which gave the minimum in eq(8.24). If any of the  $\lambda_i$  are negative then the corresponding indices are removed from the active set list.

The algorithm is

- 
1. given  $\mathbf{f}^{(1)}$  feasible and  $\mathcal{A}$  set  $k = 1$
  2. **if**  $\mathbf{h} = 0$  does not solve eq(8.20) go to 8.
  3. Compute Lagrange multipliers  $\boldsymbol{\lambda}$ .
  4. **if** all  $\lambda_i \geq 0$  terminate
  5. **else** find index  $l$  for which  $\lambda_l$  is minimum,
  6.       remove  $l$  from index set
  7. **endif**
  8. Solve eq(8.23) for  $\mathbf{h}$
  9. Find  $\tau$  from eq(8.24)
  10.  $\mathbf{f}^{(k+1)} = \mathbf{f}^{(k)} + \tau \mathbf{h}$
  11. if  $\tau < 1$  add  $p$  to  $\mathcal{A}$
  12.  $k = k + 1$ , go to 2.
-



## 9 Stochastic Optimisation

Consider an ill-posed inverse problem where a physical quantity needs to be recovered from an indirect measurement  $\mathbf{g}$  modelled by

$$\mathbf{g} = A\mathbf{f} + \epsilon, \quad (9.1)$$

where  $\mathbf{g} \in \mathbb{R}^k$  is the measurement data vector,  $\mathbf{f} \in \mathbb{R}^n$  is a discrete representation of the physical quantity of interest,  $A$  is a linear operator modelling the measurement and  $\epsilon$  is random noise taking values in  $\mathbb{R}^k$ .

In Bayesian inversion the quantities  $\mathbf{g}$  and  $\mathbf{f}$  and  $\epsilon$  are considered as random variables. The complete solution to the inverse problem is the *posterior distribution* defined using the Bayes formula:

$$\pi_{\text{post}}(\mathbf{f} | \mathbf{g}) = \frac{\pi(\mathbf{f}) \pi(\mathbf{g} | \mathbf{f})}{\pi(\mathbf{g})}, \quad (9.2)$$

where  $\pi(\mathbf{f})$  is the *prior model* representing *a priori* information about the unknown and  $\pi(\mathbf{g})$  is a normalizing constant. The conditional probability  $\pi(\mathbf{g} | \mathbf{f})$  is the *likelihood model* describing the measurement. In the case of independent white noise  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$  we have

$$\pi(\mathbf{g} | \mathbf{f}) = \pi_{\epsilon}(\mathbf{g} - A\mathbf{f}) = C \exp\left(-\frac{1}{2\sigma^2} \|A\mathbf{f} - \mathbf{g}\|_2^2\right), \quad (9.3)$$

where  $C$  is a normalization constant.

Practical measurement gives us a realization  $\hat{\mathbf{g}} \in \mathbb{R}^k$  of the random variable (9.1), which is substituted into (9.2):

$$\pi_{\text{post}}(\mathbf{f} | \hat{\mathbf{g}}) = \frac{\pi(\mathbf{f}) \pi(\hat{\mathbf{g}} | \mathbf{f})}{\pi(\hat{\mathbf{g}})}. \quad (9.4)$$

We recover  $\mathbf{f}$  from  $\hat{\mathbf{g}}$  as the conditional mean (CM) estimate

$$f^{\text{CM}} := \int_{T_n(X)} \mathbf{f} \pi_{\text{post}}(\mathbf{f} | \hat{\mathbf{g}}) d\mu(\mathbf{f}). \quad (9.5)$$

We also compute confidence intervals for uncertainty quantification.

For general references on Bayesian inversion, see [9].

### 9.1 Posterior Sampling and MCMC

Numerical evaluation of the estimate  $f^{\text{CM}}$  using (9.5) involves integration over  $n$ -dimensional space with large  $n$ . The curse of dimensionality prevents us from using integration quadratures. Instead, we approximate  $f^{\text{CM}}$  by generating a set of random samples  $\{f^{(\ell)}, \ell = 1, 2, \dots, N\}$  from the posterior density  $\pi_{\text{post}}(\mathbf{f} | \hat{\mathbf{g}})$  in such a way that

$$\frac{1}{N} \sum_{\ell=1}^N f^{(\ell)} \approx \int_{\mathbb{R}^n} \mathbf{f} \pi_{\text{post}}(\mathbf{f} | \hat{\mathbf{g}}) d\mathbf{f} = f^{\text{CM}}. \quad (9.6)$$

Confidence intervals can be approximately calculated from the variances of the sequence  $\{f^{(\ell)}, \ell = 1, 2, \dots, N\}$ .

As a typical example Markov chain Monte Carlo (MCMC) method for creating the sequence  $\{f^{(\ell)}, \ell = 1, 2, \dots, N\}$  we describe the Metropolis-Hastings algorithm. The initial state  $f^{(1)}$  is arbitrary. Assume that we know  $f^{(\ell-1)}$  and need to construct  $f^{(\ell)}$ . We create a candidate  $c \in \mathbb{R}^n$  by drawing it randomly from a so-called *proposal distribution*, for example from the Gaussian distribution  $\mathcal{N}(f^{(\ell-1)}, \beta^2)$  with a suitable standard deviation  $\beta > 0$ . Now we compare the posterior probability  $\pi_{\text{post}}(c | \widehat{\mathbf{g}})$  of the candidate to  $\pi_{\text{post}}(f^{(\ell-1)} | \widehat{\mathbf{g}})$ . If the candidate is more probable than  $f^{(\ell-1)}$ , we set  $f^{(\ell)} := c$ . Otherwise we choose one of the following two options with certain probabilities:  $f^{(\ell)} := f^{(\ell-1)}$  or  $f^{(\ell)} := c$ .

Other popular MCMC methods include the Gibbs sampler and various adaptive methods, such as DRAM and SCAM.

**References for this section :** [9].

## 10 Example Problems

### 10.1 Denoising

Consider an image corrupted by noise

$$\mathbf{g} = \mathbf{f} + \mathbf{n} \quad (10.1)$$

This takes the same general form as our linear inverse problems, with the forward mapping being the identity, if we look for a solution

$$\mathbf{f}_* = \arg \min_{\mathbf{f}} \left[ \frac{1}{2} \|\mathbf{g} - \mathbf{f}\|^2 + \frac{\alpha}{2} \|\mathbf{f}\|_{\Gamma}^2 \right] \quad (10.2)$$

The direct solution is found by solving

$$(\mathbf{I} - \alpha \Gamma) \mathbf{f}_* = \mathbf{g} \quad (10.3)$$

#### 10.1.1 Comment on Fourier Denoising

Suppose that the regulariser is a linear convolution operator, i.e.

$$\Psi(\mathbf{f}) = \|K * f\|^2 = \int_{\Omega} |K(x) * f(x)|^2 dx = \int_{\Omega} \hat{K}(k) \hat{f}(k) dk \quad (10.4)$$

Then eq.(10.3) is directly solved by

$$\hat{f}_* = \frac{\hat{g}(k)}{1 + \alpha \hat{K}^2(k)} \quad (10.5)$$

E.g. for Tikhonov 1st order regularisation

$$\hat{f}_* = \frac{\hat{g}(k)}{1 + \alpha k^2}$$

which is the classical Weiner filter.

#### 10.1.2 Series Solutions

Provided that the matrix norm  $\alpha \|\Gamma\| < 1$  we can develop a Neumann series for eq.(10.3) :

$$\mathbf{f}_* = \mathbf{g} + \alpha \Gamma \mathbf{g} + (\alpha \Gamma)^2 \mathbf{g} + \dots \quad (10.6)$$

Compare this to the Krylov subspace series

$$\mathbb{K}_{\alpha} = \{\mathbf{g}, (\mathbf{I} - \alpha \Gamma) \mathbf{g}, (\mathbf{I} - \alpha \Gamma)^2 \mathbf{g}, \dots\} \quad (10.7)$$

Clearly the space spanned by  $\mathbb{K}_\alpha$  is the same as the space spanned by

$$\mathbb{K}_{\text{Neum}} = \{\mathbf{g}, \Gamma \mathbf{g}, \Gamma^2 \mathbf{g}, \dots\} \quad (10.8)$$

### 10.1.3 Iterative Solutions

The gradient of  $\Psi(\mathbf{f})$  in eq.(10.3) is given by

$$\Psi'(\mathbf{f}) = (\mathbf{I} + \alpha \Gamma) \mathbf{f} - \mathbf{g}.$$

We can therefore form an iterative gradient descent solution by

$$\begin{aligned} \mathbf{f}^{(1)} &= \mathbf{g} \\ \mathbf{f}^{(n+1)} &= \mathbf{f}^{(n)} - \tau (\mathbf{I} + \alpha \Gamma) \mathbf{f}^{(n)} + \tau \mathbf{g} \end{aligned}$$

The step length  $\tau$  can be chosen optimally, but we can also consider this iteration to be an evolution in time:

$$\frac{\partial \mathbf{f}}{\partial t} = (\mathbf{g} - \mathbf{f}) - \alpha \Gamma \mathbf{f} \quad (10.9)$$

We note the comparison to a simple image scale space evolution scheme:

$$\frac{\partial \mathbf{f}}{\partial t} = -\Gamma \mathbf{f} \quad (10.10)$$

thus the gradient descent scheme can be considered as an image evolution scheme, followed by adding back of the residual error  $(\mathbf{g} - \mathbf{f})$ .

$$\begin{aligned} \mathbf{f}^{(n+1/2)} &= \mathbf{f}^{(n)} - \Delta t \Gamma \mathbf{f}^{(n)} \\ \mathbf{f}^{(n+1)} &= \mathbf{f}^{(n+1/2)} + (\mathbf{g} - \mathbf{f}^{(n)}) \end{aligned}$$

This is an example of a *Split Operator* method. Note that the fictitious time step  $\Delta t$  takes the place of a regularisation parameter  $\alpha$ . The time evolution concept may be particularly useful if we are considering non-linear regularisation. For example Total Variation will take the form

$$\frac{\partial \mathbf{f}}{\partial t} = \nabla \cdot \kappa(\mathbf{f}) \nabla \mathbf{f} \quad (10.11)$$

and we make take arbitrarily small steps in time to capture the non-linearity.

Finally we should note that image evolution schemes using an explicit scheme need to take account of the Courant-Freidrich-Levy (CFL) stability conditions which control stability. This sets an upper bound on the time step that can be taken. This limit is superceded if an explicit or semi-implicit evolution scheme is used.

## 10.2 The Radon Transform

The forward projection operator of the Radon transform in 2D as an integral operator mapping a function  $f(x, y)$  to a function  $g(s, \theta)$ , by integrating along the parallel lines  $\mathbf{r} \cdot \hat{\mathbf{n}} = s$ .

$$g(s, \theta) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \delta(s - (y \cos \theta - x \sin \theta)) dx dy = \int_{\mathbf{r} \cdot \hat{\mathbf{n}} = s} f(x, y) d\ell \quad (10.12)$$

The adjoint of the Radon Transform maps a function in  $(s, \theta)$  to one in  $(x, y)$ ; it is known as *Back-projection*

$$h(x, y) = \int_{-\infty}^{\infty} \int_0^{\pi} b(s, \theta) \delta(s - (y \cos \theta - x \sin \theta)) d\theta ds \quad (10.13)$$

The 1D Fourier Transform  $\mathcal{F}_{s \rightarrow k}$  of the Radon Transform of a 2D function  $f(x, y)$  is the same as the 2D Fourier Transform  $\mathcal{F}_{x \rightarrow k_x, y \rightarrow k_y}$  of  $f(x, y)$  sampled along the line  $k_x = k \cos \theta, k_y = k \sin \theta$ .

$$\int_{-\infty}^{\infty} e^{-iks} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \delta(s - (y \cos \theta - x \sin \theta)) dx dy ds = \quad (10.14)$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-iks} f(x, y) \delta(s - (y \cos \theta - x \sin \theta)) ds dx dy = \quad (10.15)$$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-ik(x \sin \theta - y \cos \theta)} f(x, y) dx dy = \hat{F}(k \sin \theta, k \cos \theta) \quad (10.16)$$

As an example, consider the Radon Transform of a 2D Gaussian function  $f(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$  which is a Gaussian function  $g(s, \theta) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{s^2}{2\sigma^2}}$  that does not depend on  $\theta$ . In this particular case  $\hat{G}(k, \theta) = \frac{\sigma}{\sqrt{2\pi}} e^{-\frac{k^2\sigma^2}{2}}$  which obviously is also the radial form of the 2D Fourier Transform  $\hat{F}(k_x, k_y) = \frac{\sigma^2}{2\pi} e^{-\frac{(k_x^2+k_y^2)\sigma^2}{2}}$

If we back project we get

$$h(x, y) = \int_0^{\pi} \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(y \cos \theta - x \sin \theta)^2}{2\sigma^2}} d\theta \quad (10.17)$$

and we clearly see that the back projection operator is not itself the inverse of the forward projection.

An exact reconstruction formula can be found by either i) back projection followed by convolution with a 2D filter or ii) convolution with a 1D filter of the Radon transform, followed by back-projection. To see this, we notice that integration of the back projection over all angles gives a weighting of  $\frac{1}{|k|}$  in the Fourier domain. We may therefore filter by the inverse Fourier transform of  $|k|$  in either 2D or 1D i)  $h(x, y) = \mathcal{F}^{-1}[1/\sqrt{(k_x^2 + k_y^2)}]$  or ii)  $h(s) = \mathcal{F}^{-1}[1/|k|]$ . These are approximately Laplacian functions (differentiating). The 2D case is hard to implement, since it is

circularly symmetric and requires interpolation. The 1D case is easier. In both cases they need to be adopted by smoothing otherwise the high frequencies will explode.

## A Further Notes on Differentiation

### A.1 Simple derivatives

Simple derivatives are examples of the general definition when  $X$  and  $Y$  are both  $\mathbb{R}$  and therefore the mapping  $A$  is a simple function  $f$ . Notice that the most general linear operator from  $\mathbb{R}$  to  $\mathbb{R}$  can be described by a simple scalar (the slope of the linear function), and the operation of  $A'(x)h$  is just scalar multiplication. Thus, the derivative should be the scalar  $f'(x)$  where

$$f(x+h) = f(x) + f'(x)h + \mathcal{O}(h).$$

From this definition, we can find this scalar  $f'(x)$  for each value of  $x$  using limits:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

This is probably a more familiar definition but you should see now how it can be thought of as a computational device for the more general definition. Notice that the key to getting a handle on the Fréchet derivative is to first understand what a general linear operator looks like when going from  $X$  to  $Y$  and then trying to figure out how to describe the derivative linear operator.

### A.2 Derivatives of vector-valued functions

Suppose now that  $X = \mathbb{R}$  and  $Y = \mathbb{R}^n$ . We usually write  $\mathbf{f} : \mathbb{R} \mapsto \mathbb{R}^n$ . A general linear mapping from  $\mathbb{R}$  to  $\mathbb{R}^n$  could be denoted with  $\mathbf{c}x$  where  $\mathbf{c}$  is a vector in  $\mathbb{R}^n$  and  $\mathbf{c}x$  denotes scalar-vector multiplication. Therefore, from the general definition we ought to be able to write

$$\mathbf{f}(x+h) = \mathbf{f}(x) + \mathbf{f}'(x)h + \mathcal{O}(h).$$

Using limits we can as before come up with

$$\mathbf{f}'(x) = \lim_{h \rightarrow 0} \frac{\mathbf{f}(x+h) - \mathbf{f}(x)}{h}.$$

which is a straightforward extension of the simple case. In practice we usually recognize that

$$\mathbf{f}'(x) = [f'_1(x), f'_2(x), \dots, f'_n(x)]^T$$

where  $f'_k(x)$  are simple derivatives so we can compute them using the familiar rules.

### A.3 Derivatives with respect to vectors (matrices)

Now, suppose the mapping is  $f : \mathbb{R}^n \mapsto \mathbb{R}$  (a multivariable function). The most general linear mapping from  $\mathbb{R}^n$  to  $\mathbb{R}$  is described by a constant length- $n$  vector,  $\mathbf{c}$ , so that

$$\mathbf{c}\mathbf{x} = \mathbf{c}^T \mathbf{x} = \langle \mathbf{x}, \mathbf{c} \rangle$$

is the value of the linear mapping. Thus, we ought to be able to define a derivative  $f'(x)$  using a length  $n$  vector. It is common in this case to use alternative notations to describe this vector.

$$f'(\mathbf{x}) \equiv \frac{\partial f}{\partial \mathbf{x}} \equiv \nabla f$$

and we sometimes call this linear operator the *gradient* of the vector. Thus, we can write

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla^T f(\mathbf{x}) \mathbf{h} + \mathcal{O}(\|\mathbf{h}\|).$$

In practice how can we express this gradient vector? We can find it projected against any direction of interest by choosing  $\mathbf{h}$  to be a scaled vector:  $\mathbf{h} = h\mathbf{q}$  where  $\mathbf{q}$  is a unit-vector in the direction of interest:

$$f(\mathbf{x} + h\mathbf{v}) = f(\mathbf{x}) + h \langle \mathbf{v}, \nabla f \rangle + \mathcal{O}(h_i)$$

and using limits we can see that

$$(\nabla f, \mathbf{v}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x})}{h}.$$

If  $\mathbf{v}$  is chosen to be one of the standard Cartesian basis vectors,  $\mathbf{e}_i$  for  $\mathbb{R}^n$  (i.e. a one in one component and zeros in the rest) then we obtain the standard definition of partial derivative

$$(\nabla f, \mathbf{e}_i) = (\nabla f)_i = \frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}.$$

Here, the computation of the partial derivative is obtained from a more general perspective which shows that  $\nabla f$  is the more general object and  $\frac{\partial f}{\partial x_i}$  is just a projection of this vector against a standard basis. Under-standing the generality will help you make sense of the derivative when we get to mappings on uncountably infinite-dimensional spaces.

**Einstein notation** Suppose we are given a formula involving a vector or a matrix and we are trying to find the derivative with respect to the vector or matrix. This is usually best handled by using summation notation (Einstein notation) and using the rules for simple differentiation. Some texts define a collection of special rules for vector (matrix) differentiation. These can be difficult to remember and keep straight. Fortunately, they are superfluous and unnecessary if you get a little comfortable with Einstein notation.

For example, consider the following formula

$$f(\mathbf{x}) = (\mathbf{y} - \mathbf{A}\mathbf{x})^T \mathbf{B}(\mathbf{y} - \mathbf{A}\mathbf{x}) + \mathbf{x}^T \mathbf{C}\mathbf{x} + \mathbf{y}^T \mathbf{D}\mathbf{y}.$$



To find  $\frac{\partial f}{\partial \mathbf{x}}$  we could try to remember the various rules associated with differentiation with respect to vectors or, alternatively, we just write the expression using summation notation

$$f(\mathbf{x}) = \sum_{i,k} \left( y_i - \sum_j A_{ij} x_j \right) B_{ik} \left( y_k - \sum_l A_{kl} x_l \right) + \sum_{i,j} x_i C_{ij} x_j + \sum_{i,j} y_i D_{ij} y_j .$$

Now, we can use the rules of simple differentiation to find ( $\delta_{ij} = 0$  unless  $i = j$  when it is 1).

$$\begin{aligned} \frac{\partial f}{\partial x_n} &= \sum_{i,k} \left( y_i - \sum_j A_{ij} x_j \right) B_{ik} - (A_{kl} \delta_{ln}) + \sum_{i,k} \left( - \sum_j A_{ij} \delta_{jn} \right) B_{ik} (y_k - A_{kl} x_l) \\ &\quad + \sum_{i,j} (x_i C_{ij} \delta_{jn} + \delta_{in} C_{ij} x_j) . \end{aligned} \quad (\text{A.1})$$

After you do this several times, you start to notice that the  $\sum$  sign is getting in the way. A simplified notation can be obtained by eliminating the sum and just remembering that repeated indices in the formula imply a sum over that index variable. Thus, the formula becomes

$$f(\mathbf{x}) = (y_i - A_{ij} x_j) B_{ik} (y_k - A_{kl} x_l) + x_i C_{ij} x_j + y_i D_{ij} y_j$$

and

$$\begin{aligned} \frac{\partial f}{\partial x_n} &= (y_i - A_{ij} x_j) B_{ik} (-A_{kl} \delta_{ln}) + (-A_{ij} \delta_{jn}) B_{ik} (y_k - A_{kl} x_l) + x_i C_{ij} \delta_{jn} + \delta_{in} C_{ij} x_j \\ &= -(y_i - A_{ij} x_j) B_{ik} A_{kn} - A_{in} B_{ik} (y_k - A_{kl} x_l) + x_i C_{in} + C_{nj} x_j \\ &= -y_i B_{ik} A_{kn} - A_{in} B_{ik} y_k + A_{ij} B_{ik} A_{kn} x_j + A_{in} B_{ik} A_{kl} x_l + x_i C_{in} + C_{nj} x_j \\ &= -A_{kn} (B_{ik} + B_{ki}) y_i + (A_{kn} (B_{ik} + B_{ki}) A_{ij} + C_{nj} + C_{jn}) x_j . \end{aligned}$$

Notice that in the last line we freely changed the indexes of summation to get the form into something we could more easily map back into matrix-vector notation. Now, if desired we can translate this back into matrix-vector notation.

$$\nabla f(\mathbf{x}) \equiv \frac{\partial f}{\partial \mathbf{x}} = -\mathbf{A}^T (\mathbf{B}^T + \mathbf{B}) \mathbf{y} + [\mathbf{A}^T (\mathbf{B}^T + \mathbf{B}) \mathbf{A} + \mathbf{C} + \mathbf{C}^T] \mathbf{x} .$$

To effectively use Einstein notation you have to get quite comfortable with the following table.

Matrix Notation	Index Notation
F	$F_{ij}$
G	$G_{jk}$
FG	$F_{ij} G_{jk}$
$F^T G$	$F_{ji} G_{jk}$
$FG^T$	$F_{ij} G_{kj}$
$F^T G^T$	$F_{ji} G_{kj}$

The first column gives a result in matrix-notation and the second column provides the equivalent in index notation. The important thing to look at is the location of the repeated index for help in understanding whether a particular sum translates to a matrix or its transpose.

## A.4 Derivatives with respect to mappings from $\mathbb{R}^m$ to $\mathbb{R}^n$

Using our general definition we can make clear sense of what is meant by the derivative of a mapping from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ . Generally, we denote this mapping by  $\mathbf{f} : \mathbb{R}^m \mapsto \mathbb{R}^n$  with  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  where  $\mathbf{y}$  is an  $n$ -dimensional vector and  $\mathbf{x}$  is an  $m$ -dimensional vector. The most general linear mapping between these spaces can be written as  $L\mathbf{x}$  where  $L$  can be described by an  $n \times m$  matrix  $L$ . For example, if  $\mathbf{x}$  is projected against an orthonormal basis  $\{\mathbf{v}_j; j = 1 \dots m\}$  of  $\mathbb{R}^m$

$$\mathbf{x} = \sum_{j=1}^m x_j \mathbf{v}_j$$

and if  $L\mathbf{x}$  is projected against an orthonormal basis  $\{\mathbf{u}_i; i = 1 \dots n\}$  of  $\mathbb{R}^n$  then

$$L\mathbf{x} = \sum_{i=1}^n y_i \mathbf{u}_i = \sum_{j=1}^m x_j L\mathbf{v}_j$$

If we take the inner product of this result with  $\{\mathbf{u}_i\}$  then we get for each  $i = 1 \dots n$

$$y_i := (L\mathbf{x}, \mathbf{u}_i) = \sum_{j=1}^m x_j (L\mathbf{v}_j, \mathbf{u}_i) = \sum_{j=1}^m L_{ij} x_j$$

This is a matrix-vector equation  $\mathbf{y} = L\mathbf{x}$  where  $L_{ij} = (\mathbf{u}_i, L\mathbf{v}_j)$ . Very often we choose  $\mathbf{v}_j$  and  $\mathbf{u}_i$  to be the standard Euclidean bases and simply describe  $L$  using its standard matrix representation.

Thus, according to the general definition we should be able to find a derivative operator  $D_{\mathbf{x}} \equiv \mathbf{f}'(\mathbf{x})$  such that

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x}) + D_{\mathbf{x}}\mathbf{h} + \mathcal{O}(\|\mathbf{h}\|).$$

We can describe this linear operator  $D_{\mathbf{x}}$  using  $n \times m$  scalars. Most commonly, the standard Euclidean bases are chosen to represent both  $\mathbf{x}$  and  $\mathbf{y}$  and so the matrix entries can be expressed in terms of specific partial derivatives. Deriving a specific form for these matrix components using limits is left as an exercise for the reader.

There are two commonly used multivariable derivatives that are related to the general Fréchet derivative. These are explained in the next two sections. These two derivatives are quite useful to help explain Newtonian physics but from an optimization standpoint they are not particularly useful and so we don't deal with them. They are included here only to show you how they fit in to the overall structure we are building.

## Divergence

For the special case of a mapping from  $\mathbb{R}^n$  to  $\mathbb{R}^n$  it is sometimes convenient (particularly for Physics applications) to define the divergence of the mapping

$$\nabla \cdot \mathbf{f} = \text{Tr}(\mathbf{D}_{\mathbf{x}}),$$

where

$$\text{Tr}(\mathbf{A}) = \sum_j (\mathbf{u}_j, \mathbf{A}\mathbf{u}_j)$$

where  $\mathbf{u}_j$  is any orthonormal basis. For a standard Cartesian bases this is

$$\nabla \cdot \mathbf{f} = \sum_{i=1}^n \frac{\partial f_i}{\partial x_i} = f_{i,i} = D_{ii}(\mathbf{x}),$$

where  $f_{i,j} = \frac{\partial f_i}{\partial x_j}$  is another part of Einstein notation along with the implied sum over  $i$ . Notice also that we have used the Fréchet derivative matrix:  $\{\mathbf{D}_{\mathbf{x}}\}_{ij} = D_{ij}(\mathbf{x})$ .

## Curl

Another commonly used derivative that is applicable only for mappings from  $\mathbb{R}^3$  to  $\mathbb{R}^3$  is the curl. In Einstein notation in Cartesian coordinates the curl of  $\mathbf{f}$  is a vector  $\mathbf{g}$  defined as

$$\begin{aligned} \mathbf{g} &= \nabla \times \mathbf{f} \\ g_i &= \epsilon_{ijk} \frac{\partial f_k}{\partial x_j} = \epsilon_{ijk} f_{k,j} \\ \mathbf{g} &= \boldsymbol{\epsilon} : \mathbf{D}_{\mathbf{x}} \end{aligned}$$

where there is an implied sum over both  $j$  and  $k$ ,  $f_{k,j} = \frac{\partial f_k}{\partial x_j} = D_{kj}(\mathbf{x})$  and  $\epsilon_{ijk}$  is the  $3 \times 3 \times 3$  Levi-Civita tensor which is defined to be 0 if any of the  $i, j$ , or  $k$  are repeated indices, 1 if  $(i, j, k)$  is an even permutation of (1, 2, 3) and -1 if  $(i, j, k)$  is an odd permutation of (1, 2, 3). For example,  $\epsilon_{113} = 0$  and  $\epsilon_{231} = 1$  but  $\epsilon_{132} = -1$ . In particular,  $\epsilon_{123} = \epsilon_{231} = \epsilon_{312} = 1$  while  $\epsilon_{321} = \epsilon_{213} = \epsilon_{132} = -1$  and the rest of the components of  $\epsilon_{ijk}$  are 0.

The notation  $\mathbf{A} : \mathbf{B}$  to denote a tensor-matrix multiplication that involves two sums. This kind of notation gets bulky (thus the use of  $\nabla \times \mathbf{f}$  in vector notation).

## B Variational Methods

### B.1 Generalised Differentiation

The reason derivatives arise in practice is that all derivatives try to capture the idea of an “instantaneous” rate-of-change of a mapping with respect to its input. In other words, the derivative of a mapping is another mapping that shows how the original mapping changes when the input change. The word “instantaneous” is precisely defined using the concept of limits. If we have some mapping and we are trying to find some optimum value for that mapping then at the optimum value, the instantaneous rate of change of the mapping will be zero. Thus, finding the zeros of the derivative of the mapping will be a method of finding optima. This procedure is the basic principle of calculus for mappings from real-numbers to real-numbers (regular functions), but also extends to all kinds of mappings such as those based on vectors, functions (variational calculus), complex-valued variables, and even operators.

#### B.1.1 Definitions

Here we give the general definition of a (Fréchet) derivative of an arbitrary mapping from a normed vector space  $X$  to a normed vector space  $Y$ . The next sections will try to make better sense of this definition for the more common mappings we encounter. Suppose we have the mapping  $A : X \mapsto Y$ . The derivative of this operator at  $x$  is denoted  $A'(x)$  and is defined as the linear (and continuous, thus bounded) operator from  $X$  to  $Y$  such that

$$A(x + h) = A(x) + A'(x)h + \mathcal{O}(\|h\|) \quad (\text{B.1})$$

where by definition of little-o notation:<sup>2</sup>

$$\lim_{\|h\| \rightarrow 0} \frac{\mathcal{O}(\|h\|)}{\|h\|} = 0 \quad (\text{B.2})$$

(note that this limiting value of 0 is the 0-vector of the space  $Y$ ). The notation  $A'(x)h$  means operate on  $h$  with the (linear) derivative operator  $A'(x)$ . Exactly what this operation is depends on the specifics of  $A$  and the spaces  $X$  and  $Y$ . For this general definition  $x$  and  $h$  are elements of  $X$  and could be for example scalars, vectors (matrices), or functions. In other words,  $x$  and  $h$  could be one-dimensional, countably-dimensional, or uncountably-dimensional. Also, the addition on the left-hand side of the definition occurs in  $X$  while the addition in the right-hand side of the definition occurs in  $Y$ .

Technically this definition defines the right-derivative and we should also define the left-derivative using a backward difference and then say that the derivative exists only if these two values are the same. Normally we will deal with smooth mappings and thereby avoid this technicality.

---

<sup>2</sup>Little-o notation is useful as a container for arbitrarily complex expressions which vanish at a limit. For our purposes,  $\mathcal{O}(g)$  is such that  $\mathcal{O}(g)/g \rightarrow 0$  when  $g \rightarrow 0$ .

Notice what happens if  $A$  is a linear operator. Then, we already know that

$$A(x + h) = Ax + Ah$$

and thus we see that  $A'(x) = A$  for all  $x$ .

### B.1.2 Mappings from $X$ to $\mathbb{R}$

Suppose  $X$  is an infinite-dimensional space (like the space of square-integrable functions) and the mapping is from  $X$  to the real numbers. This is what we have called a functional, but it lives happily in our general notion of a mapping from one normed vector space to another. Finding a derivative of a functional is sometimes called variational calculus, but again you can see that it is just an example of a single definition of derivative.

The general derivative definition says we should be able to define the (Fréchet) derivative of this functional with

$$\Psi(f + h) = \Psi(f) + \Psi'(f)h + \mathcal{O}(\|h\|),$$

where  $\Psi'(f)$  is some linear operator. How can we describe all the linear operators that map from  $X$  to  $\mathbb{R}$ . The simplest way is to choose a basis of  $X$  and tabulate what happens to the basis vectors through the linear operator. So, for example, write

$$f(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) d\mathbf{y}$$

and we can see that the linear operator  $P$  will map  $f$  to the number

$$\begin{aligned} q &= Pf = \int P\delta(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) d\mathbf{y} \\ &= \int p(\mathbf{y}) f(\mathbf{y}) d\mathbf{y} = \langle p, f \rangle, \end{aligned}$$

where  $p(\mathbf{y}) = P\delta(\mathbf{x} - \mathbf{y})$  is the number that  $P$  maps  $\delta(\mathbf{x} - \mathbf{y})$  to, for a given  $\mathbf{y}$ . Thus, any linear operator,  $P$  from  $X$  to  $\mathbb{R}$  can be described using a function  $p(\mathbf{y})$ . As a result we could represent the derivative operator  $\Psi'(f)$  using a function  $\delta\Psi_f(\mathbf{y})$  with

$$\Psi'(f)h \equiv \int \delta\Psi_f(\mathbf{y}) h(\mathbf{y}) d\mathbf{y}.$$

Therefore,

$$\Psi(f + h) = \Psi(f) + \int \delta\Psi_f(\mathbf{y}) h(\mathbf{y}) d\mathbf{y} + \mathcal{O}(\|h\|).$$

To find  $\delta\Psi_f(\mathbf{y})$  we choose  $h(\mathbf{y}) = \epsilon\delta(\mathbf{x} - \mathbf{y})$  and by taking limits note that

$$\begin{aligned} \delta\Psi_f(\mathbf{y}) &= \lim_{\epsilon \rightarrow 0} \frac{\Psi(f(\mathbf{x}) + \epsilon\delta(\mathbf{x} - \mathbf{y})) - \Psi(f(\mathbf{x}))}{\epsilon} \\ &\equiv \frac{\partial\Psi}{\partial f}. \end{aligned}$$

Notice that  $\Psi'(f)$  represents the operator while  $\delta\Psi_f(\mathbf{y})$  or simply  $\delta\Psi$  represents the function that implements the operator so that

$$\Psi'(f)h = \langle \delta\Psi, h \rangle .$$

In calculations it is especially important to distinguish these notations.

### B.1.3 General Mappings

We have now shown how the definition of the (Fréchet) derivative of a mapping  $A : X \mapsto Y$  can be defined as the linear operator satisfying the equation:

$$A(f + h) = A(f) + A'(f)h + \mathcal{O}(\|h\|) \quad \text{as } \|h\| \rightarrow 0 .$$

We have seen how this general definition specializes for the more familiar mappings we are accustomed to and how the Fréchet derivative corresponds to our more familiar notions of derivative. We are now ready to get our minds around  $A'(f)$  in the most general sense.

The first question we ask is how do we describe a general linear operator mapping  $X$  to  $Y$ . We already have a couple of representations. The first representation uses the kernel  $K(\mathbf{x}; \mathbf{x}')$  so that the linear operator  $P$  is described as

$$(Pf)(\mathbf{x}) = \int K(\mathbf{x}; \mathbf{x}')f(\mathbf{x}')d\mathbf{x}' \quad (\text{B.3})$$

where  $K(\mathbf{x}; \mathbf{x}')$  is the response of the operator to an impulse  $\delta(\mathbf{x} - \mathbf{x}_0)$ . We also have the spectral representation where we note that

$$K(\mathbf{x}, \mathbf{x}') = \sum_i w_i \mathbf{u}_i(\mathbf{x}) \mathbf{v}_i(\mathbf{x}')$$

so that

$$(Pf)(\mathbf{x}) = \sum_i w_i \langle f, \mathbf{v}_i \rangle \mathbf{u}_i(\mathbf{x}) . \quad (\text{B.4})$$

Thus, describing a general linear operator amounts to finding the kernel  $K(\mathbf{x}; \mathbf{x}')$  and/or its spectral representation. Consequently, we could write the defining equation for the (Fréchet) derivative as

$$[A(f + h)](\mathbf{y}) = [A(f)](\mathbf{y}) + \int \delta A(\mathbf{y}; \mathbf{x}')h(\mathbf{x}')d\mathbf{x}' + \mathcal{O}(\|h\|) ,$$

where  $\delta A$  is the kernel of the linear operator  $A'(f)$ . By choosing  $h(\mathbf{x}) = \epsilon\delta(\mathbf{x} - \mathbf{z})$  and taking limits we could find

$$\delta A(\mathbf{y}; \mathbf{z}) = \lim_{\epsilon \rightarrow 0} \frac{[A(f(\mathbf{x}) + \epsilon\delta(\mathbf{x} - \mathbf{z}))](\mathbf{y}) - [A(f(\mathbf{x}))](\mathbf{y})}{\epsilon} .$$

The Fréchet derivative is important largely for its role in optimization theory, but it also plays a role in inverse problems in general as it shows us the local linear behavior of an arbitrary non-linear

map. Consequently, the local SVD of the Fréchet derivative gives us a clear picture of the amount of information loss occurring for a particular input.

#### B.1.4 Higher Order Derivatives

Higher order (Fréchet) derivatives are defined recursively so that  $A''(x)$  is the linear operator mapping from  $X$  to  $Z$  (where  $Z$  is the space of all linear operators from  $\mathbf{X}$  to  $\mathbf{Y}$ ) such that

$$A'(x+k) = A'(x) + A''(x)k + \mathcal{O}(\|h\|).$$

You can also think about  $A''(x)$  as defining a bilinear map from  $X \times X$  to  $Y$  so that  $(h, k) \mapsto (A''(x)k)h$ . It can be shown that this mapping is symmetric:

$$(A''(x)k)h = (A''(x)h)k.$$

We can continue this game and define a countably infinite number of higher-order Fréchet derivatives. For fun, try to come up with a general Taylor-series expansion of an arbitrary mapping in terms of these higher-order Fréchet derivatives. To be specific, try to show that we could represent an arbitrary mapping as

$$A(x+h) = A(x) + A'(x)h + c_2(A''(x)h)h + c_3((A'''(x)h)h)h + \dots,$$

where  $c_k = \frac{1}{k!}$ .

#### B.1.5 Calculus of Fréchet Derivatives

**Example Rules** This definition as a limit is quite similar to the simple limit definition for the derivative. This means that all the familiar rules of calculus apply. A maximum or minimum (or in action point) occurs when  $\delta\Psi_f(\mathbf{y}) = 0$  and we can check between these by looking at  $\Psi''$ .

Other rules that are useful:

- Chain Rule: Suppose a new functional is  $\Psi_1 = g(\Psi(f))$  then

$$\delta\Psi_1 = \delta g(\Psi(f)) = g'(\Psi(f))\delta\Psi$$

- Quotient rule

$$\delta \left[ \frac{\Psi_1}{\Psi_2} \right] = \frac{\Psi_2 \delta\Psi_1 - \Psi_1 \delta\Psi_2}{\Psi_2^2}$$

- Product rule

$$\delta[\Psi_1 \Psi_2] = \delta\Psi_1 \Psi_2 + \Psi_1 \delta\Psi_2$$

- Sum rule

$$\delta[a\Psi_1 + b\Psi_2] = a\delta\Psi_1 + b\delta\Psi_2$$

**Specific Derivatives** The derivative of a constant is zero

$$\delta a = 0.$$

In addition when  $\Psi(f) = f(y)$  for a particular value of  $y$  then

$$\delta\Psi = \delta(y).$$

When  $\Psi(f) = \int f(x)g(x)dx = \langle f, g \rangle$  we get

$$\begin{aligned}\delta\Psi &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left[ \int (f(x)g(x) + \epsilon\delta(x-y)g(x)) dx - \int f(x)g(x)dx \right] \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \int \epsilon\delta(x-y)g(x)dx \\ &= g(y).\end{aligned}$$

Similarly, when

$$\Psi(f) = \int A(f(x))B(f(x))dx = \langle A(f), B(f) \rangle \quad (\text{B.5})$$

we get

$$\begin{aligned}\delta\Psi &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\langle A(f + \epsilon\delta), B(f + \epsilon\delta) \rangle - \langle A(f), B(f) \rangle] \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\langle A(f) + \epsilon A'(f)\delta, B(f) + \epsilon B'(f)\delta \rangle - \langle A(f), B(f) \rangle] \\ &= \langle A'(f)\delta, B(f) \rangle + \langle A(f), B'(f)\delta \rangle \\ &= \langle \delta, [A'(f)]^* B(f) \rangle + \langle [B'(f)]^* A(f), \delta \rangle \\ &= [A'(f)]^* B(f) + [B'(f)]^* A(f).\end{aligned} \quad (\text{B.6})$$

Finally, when

$$\Psi(f) = \int g[f(x)]dx$$

then

$$\begin{aligned}\delta\Psi(y) &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left[ \int g(f(x) + \epsilon\delta(x-y))dx - \int g(f(x))dx \right] \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left\{ \int [g(f(x)) + \epsilon\delta(x-y)g'(f(x)) + \mathcal{O}(\epsilon) - g(f(x))] dx \right\} \\ &= \int \delta(x-y)g'(f(x))dx \\ &= g'[f(y)].\end{aligned}$$

Using just these rules quite a few calculations can be done.



**Example : Entropy** The Shannon entropy of a function is defined

$$\Psi_{\text{Ent}}(f) = \int (f \log f - f) dx$$

using the above rules, e.g. Eq. B.6 we get

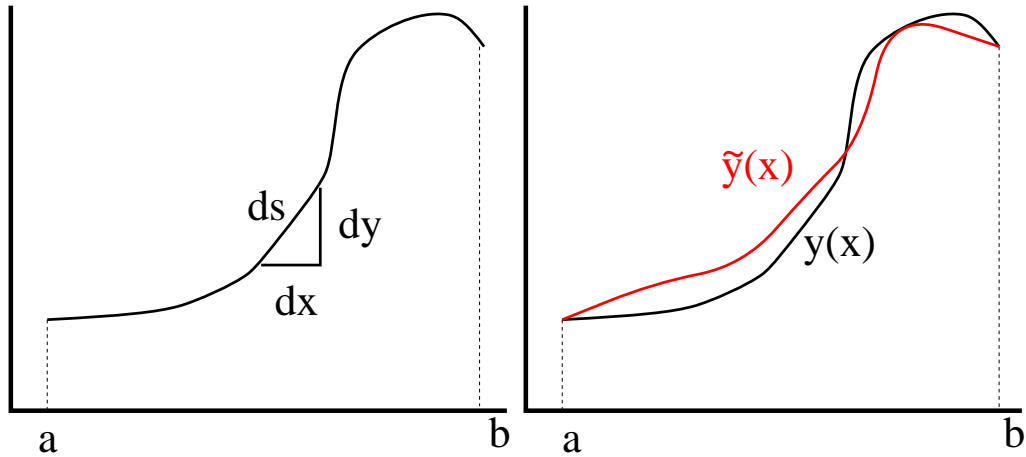
$$\delta \Psi_{\text{Ent},f} = \log f$$

So that the Fréchet derivative operator is given by

$$\Psi'(f)h = \int h \log f dx = \langle h, \log f \rangle$$

## B.2 Variational Forms and Euler-lagrange Equations

Whereas optimisation of functions (e.g. norms) over finite vector spaces may be quite familiar, it is also possible to consider infinite dimensional spaces of functions, and to look for the minimisation (or maximisation) of *functionals*.



Let us look for the stationary values (minimum, maximum, points of inflection) of the following functional

$$\Psi = \int_a^b \psi(y, y', x) dx \quad (\text{B.7})$$

Consider a perturbation from a given function  $y(x)$  in an arbitrary “direction”  $h(x)$  having the same boundary conditions

$$\tilde{y}(x) := y(x) + \epsilon h(x)$$

where  $h(a) = h(b) = 0$ . If  $\Psi$  is stationary, we define

$$\Psi(\epsilon) = \int_a^b \psi(\tilde{y}, \tilde{y}', x) dx \quad (\text{B.8})$$

and we require

$$\left. \frac{d\Psi}{d\epsilon} \right|_{\epsilon=0} = 0 \quad \forall h(x). \quad (\text{B.9})$$

If we write  $\Psi(\epsilon)$  as a Taylor series we have

$$\begin{aligned} \Psi(\epsilon) &= \int_a^b \psi(y + \epsilon h, y'(x) + \epsilon h', x) dx \\ &= \Psi(0) + \epsilon \int_a^b \left[ \frac{\partial \psi}{\partial y} h + \frac{\partial \psi}{\partial y'} h' \right] dx + \mathcal{O}(\epsilon^2) \end{aligned}$$

The requirement for stationarity is therefore

$$\Psi'(y)h = \int_a^b \left[ \frac{\partial \psi}{\partial y} h + \frac{\partial \psi}{\partial y'} h' \right] dx = 0.$$

By integration by parts and using  $h(a) = h(b) = 0$  we find

$$\Psi'(y)h = \underbrace{\left[ \frac{\partial \psi}{\partial y'} h \right]_a^b}_{=0} + \int_a^b \underbrace{\left[ \frac{\partial \psi}{\partial y} - \frac{d}{dx} \frac{\partial \psi}{\partial y'} \right]}_{=\delta\Psi} h dx = \langle \delta\Psi, h \rangle = 0$$

and since this must be true for any  $h$  we derive the *Euler-Lagrange* equation :

$$\delta\Psi := \frac{\partial \psi}{\partial y} - \frac{d}{dx} \frac{\partial \psi}{\partial y'} = 0 \quad (\text{B.10})$$

An expression such as *Eq. B.7* is known as a *Variational Form*, or sometimes as an *Energy* term because of its relation to principles such as Free Energy in classical thermodynamics. The Euler-Lagrange equation establishes an important principal in that optimisation of a variational form is always equivalent to solving a partial differential equation (PDE). A useful concept that we will see later, is that the solution to some (not all) PDEs can be made equivalent to optimising a variational form.

An alternative expression is sometimes useful :

$$\frac{\partial \psi}{\partial x} - \frac{d}{dx} \left( \psi - y' \frac{\partial \psi}{\partial y'} \right) = 0. \quad (\text{B.11})$$

To see this, expand all the derivatives in *Eq. B.11* explicitly,

$$\begin{aligned} \frac{\partial \psi}{\partial x} - \frac{d}{dx} \left( \psi - y' \frac{\partial \psi}{\partial y'} \right) &= \frac{\partial \psi}{\partial x} - \underbrace{\left( \frac{\partial \psi}{\partial x} + y' \frac{\partial \psi}{\partial y} + y'' \frac{\partial \psi}{\partial y'} \right)}_{\frac{d\psi}{dx}} + \underbrace{\left( y'' \frac{\partial \psi}{\partial y'} + y' \frac{d}{dx} \frac{\partial \psi}{\partial y'} \right)}_{\frac{d}{dx} y' \frac{\partial \psi}{\partial y'}} \\ &= -y' \left( \frac{\partial \psi}{\partial y} - \frac{d}{dx} \frac{\partial \psi}{\partial y'} \right) = 0 \end{aligned}$$

**Example 1 : Shortest path.** What is the shortest path between two points? I.e. if we consider the path as a function  $y(x)$  how do we find the one with minimum length ?

$$L = \int_a^b ds = \int_a^b (dx^2 + dy^2)^{1/2} = \int_a^b \left(1 + \left(\frac{dy}{dx}\right)^2\right)^{1/2} dx \quad (\text{B.12})$$

Obviously it is a straight line, but how do we actually prove that ? This leads us to the topic of *calculus of variations*.

For the example above, we have

$$\begin{aligned} \psi(y, y', x) &= (1 + y'^2)^{1/2} \\ \Rightarrow \frac{\partial \psi}{\partial y'} &= \frac{y'}{(1 + y'^2)^{1/2}} \end{aligned}$$

From the Euler-Lagrange equation we have

$$\begin{aligned} \frac{d}{dx} \frac{\partial \psi}{\partial y'} &= 0 \\ \Rightarrow \frac{\partial \psi}{\partial y'} &= \frac{y'}{(1 + y'^2)^{1/2}} = \text{constant} = A \end{aligned}$$

Rearranging,

$$y' = \frac{A}{(1 - A^2)^{1/2}} = m \Rightarrow y = mx + c$$

and the values of the slope  $m$  and offset  $c$  are found from the boundary conditions. [Note also that  $m = \frac{A}{(1 - A^2)^{1/2}} \Rightarrow A = \cos \theta$ ].

**Example 2 : Light path in variable refractive index medium.** Fermat's principle says that light follows the path that minimises the *time* spent in the medium

$$T = \frac{1}{c_0} \int_a^b n(x, y) (1 + y'^2)^{1/2} dx \quad (\text{B.13})$$

where  $c_0$  is the speed of light in a vacuum. In this case we make use of the second Euler-Lagrange

form, Eq. B.11 to get

$$\begin{aligned}
\frac{d}{dx} \left( \psi - y' \frac{\partial \psi}{\partial y'} \right) &= 0 \\
\Rightarrow \left( \psi - y' \frac{\partial \psi}{\partial y'} \right) &= \text{constant} = K \\
\Rightarrow n(x, y)(1 + y'^2)^{1/2} - \frac{n(x, y)y'^2}{(1 + y'^2)^{1/2}} &= K \\
\Rightarrow \frac{n(x, y)}{(1 + y'^2)^{1/2}} &= K
\end{aligned}$$

For example if  $n = e^y$  and  $a = -1$ ,  $b = 1$ , then we can find the explicit solution

$$y(x) = \frac{1}{2} \log [1 + \tan(x)^2] + \text{constant} \quad (\text{B.14})$$

where the constant is chosen to satisfy the boundary conditions

### B.3 Generalisations

We may generalise the basic variational method in several ways

**Higher order derivatives .** Suppose that

$$\Psi(y) = \int_a^b \psi(x, y, y', y'', \dots y^{(n)}) dx. \quad (\text{B.15})$$

We obtain

$$\frac{\delta \psi}{\delta y} = \frac{\partial \psi}{\partial y} - \frac{d}{dx} \frac{\partial \psi}{\partial y'} + \frac{d^2}{dx^2} \frac{\partial \psi}{\partial y''} \dots + (-1)^n \frac{d^n}{dx^n} \frac{\partial \psi}{\partial y^{(n)}} = 0 \quad (\text{B.16})$$

**More than one dependent variable.** Suppose that the function sought is a vector function of one variable (e.g. a curve in space)

$$\Psi(y_1, y_2, \dots y_n) = \int_a^b \psi(x, y_1, y'_1, y_2, y'_2, \dots y_n, y'_n) dx = \int_a^b \psi(x, \mathbf{y}, \mathbf{y}') dx \quad (\text{B.17})$$

with specified boundary conditions  $\mathbf{y}(a) = \boldsymbol{\alpha}$ ,  $\mathbf{y}(b) = \boldsymbol{\beta}$ , then the above procedure leads to the Euler-Lagrange equations

$$\frac{\delta \psi}{\delta y_i} = \frac{\partial \psi}{\partial y_i} - \frac{d}{dx} \frac{\partial \psi}{\partial y'_i} = 0 \quad (\text{B.18})$$

**More than one independent variable.** Now consider a scalar function of several variables (e.g. a surface S)

$$\Psi(y) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_m}^{b_m} \psi(x_1, x_2, \dots, x_m, y, y_{x_1}, y_{x_2}, \dots, y_{x_m}) dx_1 dx_2 \dots dx_m = \int_S \psi \left( \mathbf{x}, y, \frac{\partial y}{\partial \mathbf{x}} \right) d^m \mathbf{x} \quad (\text{B.19})$$

with specified boundary conditions  $y(\mathbf{x})|_S = h$ . The Euler-Lagrange equations are

$$\frac{\delta \psi}{\delta y} = \frac{\partial \psi}{\partial y} - \sum_{j=1}^m \frac{d}{dx_j} \frac{\partial \psi}{\partial y_{x_j}} = 0 \quad (\text{B.20})$$

**Multivariate functions of multiple variables.** Finally, consider a vector function of a vector of variables

$$\Psi(y_1, y_2, \dots, y_n) = \int_S \psi \left( \mathbf{x}, \mathbf{y}, \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right) d^m \mathbf{x} \quad (\text{B.21})$$

with specified boundary conditions  $y(\mathbf{x})|_S = \mathbf{h}$ . The Euler-Lagrange equations are

$$\frac{\delta \psi}{\delta y_i} = \frac{\partial \psi}{\partial y_i} - \sum_{j=1}^m \frac{d}{dx_j} \frac{\partial \psi}{\partial y_{i,x_j}} = 0 \quad (\text{B.22})$$

## References

- [1] J. E. Dennis and R. B. Schnabel. *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Prentice-Hall, Englewood Cliffs, New Jersey, 1983.
- [2] R. Fletcher. *Practical Methods of Optimization*. Wiley, Chichester, second edition, 1987.
- [3] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Verlag, New York, 1999.
- [4] S. Boyd and L. Vandenberghe. *Convex Optimisation*. Cambridge University Press, 2004.
- [5] C. R. Vogel. *Computational Methods for Inverse Problems*. SIAM, 2002.
- [6] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. *Numerical Recipes in C*. Cambridge University, Cambridge, England, 2nd edition, 1992.
- [7] L. N. Trefethen and D. Bau III. *Numerical Linear Algebra*. SIAM, Philadelphia, 1997.
- [8] M. Bertero and M. Boccacci. *Introduction to Inverse Problems in Imaging*. IOP Publishing Ltd, 1998.
- [9] J. Kaipio and E. Somersalo. *Statistical and Computational Inverse Problems*. Springer, New York, 2005.
- [10] G. Sapiro. *Geometric Partial Differential Equations and Image Analysis*. Cambridge Univ. Press, New York, 2001.
- [11] O. Scherzer, M. Grasmair, H. Grossauer, M. Haltmeier, and F. Lenzen. *Variational Methods in Imaging*. Springer, 2009.