# 8

## *Introduction to Retrieval Theory*

### 8.1  Measurement

A measurement describes a set of operations which determine the value of a quantity. The quantity may be a scalar, e.g. the temperature of a body, or a vector, e.g. an atmospheric temperature profile. Following *BIPM* [2008] it is helpful to define *measurand* as *the particular quantity subject to measurement* so that the phrases 'true value of a quantity' and 'value of the measurand' are synonymous.

a)
**Direct Measurement**

measurand ⟶ measurement(s) = estimate of the measurand

x?

b)
**Indirect Measurement**     **Retrieval Theory**

measurand ⟶ measurement(s) ⟶ estimate of the measurand
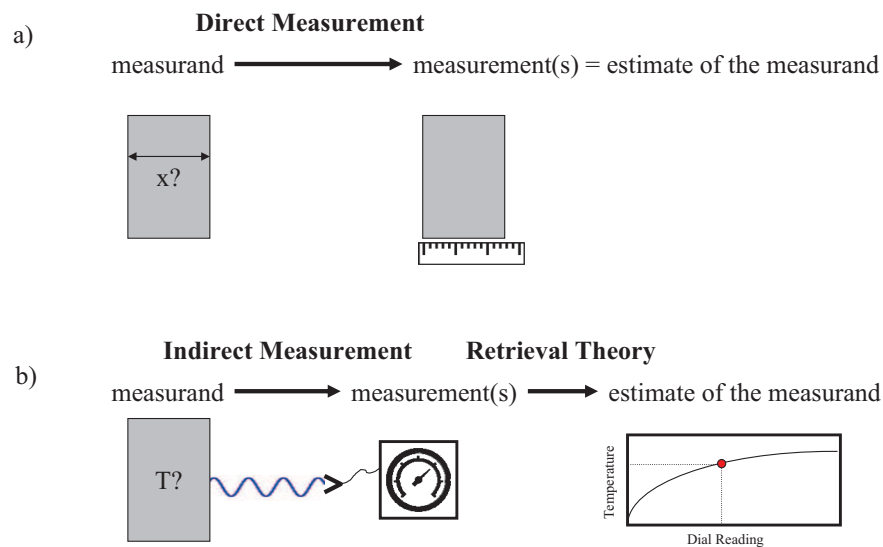
T?

Temperature

Dial Reading

**FIGURE 8.1**

The estimates of the measurand are made (a) by direct measurement or (b) through applying a retrieval process to indirect measurements. An example of a direct measurement would be the estimation of the width of this page using a ruler. An example of an indirect measurement would be estimating the temperature of this page from measurements of the radiation emitted at a specified wavelength.

Measurement includes the step of either direct or indirect measurement. If the quantity is measured without an intervening medium or process it is called a direct measurement. Very few instruments directly measure the quantity of interest: instead, most instruments make an indirect or remote measurement of some effect caused by the measurand. Inverse or retrieval methods are then used to determine the measurand from the measurements. The distinction between measurements which include either a direct or an indirect measurement step is shown schematically in Figure 8.1.

Retrieval problems have been addressed by many fields and have formed the base of several atmospheric science texts. Here, there is no intention of reviewing the field exhaustively: instead, the notation and general approach of Rodgers [*Rodgers*, 1976, 2000] is followed. This methodology is widely adopted and is a starting point from which other methods could be explored.

### 8.1.1   Mean, Variance, Covariance & the Gaussian Distribution

For a set of $N$ numbers $\{x_1, x_2, \ldots, x_N\}$ the mean value, $\mu$ or $\langle x \rangle$ is

$$\mu = \langle x \rangle = \frac{1}{N} \sum_{i=1}^{N} x_i \tag{8.1}$$

and a measure of the spread of the values is the variance, $\sigma^2$, defined by

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^{N} x_i^2 - \left( \frac{1}{N} \sum_{i=1}^{N} x_i \right)^2 = \langle x^2 \rangle - \langle x \rangle^2 \tag{8.2}$$

which is the mean square value minus the squared mean value.

For a set of $M$ vectors $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_M\}$ of length $N$ the mean value, $\langle \mathbf{x} \rangle$ is

$$\langle \mathbf{x} \rangle = \frac{1}{M} \sum_{j=1}^{M} \mathbf{x}_j \tag{8.3}$$

The mean of the $i^{\text{th}}$ component of $\mathbf{x}$ is shown as $\langle x_i \rangle$ or $\mu_i$ defined by

$$\mu_i = \langle x_i \rangle = \frac{1}{N} \sum_{j=1}^{N} x_{i,j} \tag{8.4}$$

where $x_{i,j}$ is the $i^{\text{th}}$ component of the $j^{\text{th}}$ vector.

The covariance between two variables, $x_a$,$x_b$ say, in the vector is

$$\text{Cov}(x_a, x_b) = \frac{1}{M} \sum_{k=1}^{M} (x_{a,k} - \mu_a)(x_{b,k} - \mu_b) = \langle x_a x_b \rangle - \mu_a \mu_b \tag{8.5}$$

The covariance matrix, $\mathbf{S}$, is defined by

$$S_{ab} = \text{Cov}(x_a, x_b) \quad \text{or} \quad \mathbf{S} = \begin{pmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_2, x_1) & \dots & \text{Cov}(x_N, x_1) \\ \text{Cov}(x_1, x_2) & \text{Cov}(x_2, x_2) & \dots & \text{Cov}(x_N, x_2) \\ \vdots & \vdots & & \vdots \\ \text{Cov}(x_1, x_N) & \text{Cov}(x_2, x_N) & \dots & \text{Cov}(x_N, x_N) \end{pmatrix} \quad (8.6)$$

As $\text{Cov}(x_a, x_b) = \text{Cov}(x_b, x_a)$ a covariance matrix is symmetric. It is also positive definite (REF). The diagonal elements of the covariance matrix are the variances on the measurements. The off-diagonal elements allow for correlation between the values of $\mathbf{x}$.

A Gaussian probability density function, $P(x)$, is defined

$$P(x) = N(\mu, \sigma) \quad (8.7)$$

such that $P(x)dx$ represents the probability that $x$ lies in the interval $(x, x + dx)$. An advantage of using Gaussian representations of pdfs is the algebraic convenience. For example, the mean and variance for a continuous distribution $P(x)$ are

$$\mu = \int x P(x)\,dx \quad (8.8)$$

$$\sigma^2 = \int (x - \mu)^2 P(x)\,dx \quad (8.9)$$

The product of a Gaussian with mean and standard deviation $\mu_u, \sigma_u$ with a second Gaussian with mean and standard deviation $\mu_v, \sigma_v$ is a third Gaussian, i.e.

$$\begin{aligned} P(x) &\propto \exp\left[-\frac{(x - \mu_u)^2}{2\sigma_u^2}\right] \exp\left[-\frac{(x - \mu_v)^2}{2\sigma_v^2}\right] \\ &= \exp\left[-\frac{(x - \mu_u)^2}{2\sigma_u^2} - \frac{(x - \mu_v)^2}{2\sigma_v^2}\right] \\ &= \exp\left[-\left(x - \frac{x_u/\sigma_u^2 + x_v/\sigma_v^2}{1/\sigma_u^2 + 1/\sigma_v^2}\right)^2 \frac{\sigma_u^2 + \sigma_v^2}{2\sigma_u^2\sigma_v^2}\right] \end{aligned} \quad (8.10)$$

which is more simply expressed in terms of the Gaussian parameters so that the product is a Gaussian with a mean given by

$$\mu = \frac{\mu_u/\sigma_u^2 + \mu_v/\sigma_v^2}{1/\sigma_u^2 + 1/\sigma_v^2} \quad (8.11)$$

and a variance given by

$$\frac{1}{\sigma^2} = \frac{1}{\sigma_u^2} + \frac{1}{\sigma_v^2}. \quad (8.12)$$

In $n$ dimensions the Gaussian measurement probability density function is

$$P(\mathbf{x}) = N(\langle \mathbf{x} \rangle, \mathbf{S}_x) = \frac{1}{\sqrt{(2\pi)^n}|\mathbf{S}_x|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \langle \mathbf{x} \rangle)^{\text{T}} \mathbf{S}_x^{-1}(\mathbf{x} - \langle \mathbf{x} \rangle)\right] \quad (8.13)$$

where $\mathbf{S}_x$ is the covariance matrix of $\mathbf{x}$ and $\langle \mathbf{x} \rangle$ is the vector mean. It follows that

$$\langle \mathbf{x} \rangle = \int \mathbf{x} P(\mathbf{x}) \, d\mathbf{x}. \tag{8.14}$$

### 8.1.2  Measurement Error and Uncertainty

The process of measurement is inexact, so the difference between a measured value and the measurand is called the error. Except in a few cases, the "true" value of the error is not known and its magnitude is hypothetical. There are many reasons why a measurement is uncertain. For example, error components in satellite remote sensing may include terms such as

- instrument noise,

- error arising from simplifications in radiative transfer,

- calibration error,

- error arising from the uncertainty in parameters used to derive the measurement.

A single measurement can be considered as the outcome from a trial that has sampled an infinite set of possible measurements. If we were able to perform exactly the same measurement repeatedly then a histogram of the values would tend to the measurement probability density function, $P(\mathbf{y})$, where the value of $P(\mathbf{y}) d\mathbf{y}$ is the probability that the measurement was in the multidimensional interval $(\mathbf{y}, \mathbf{y} + d\mathbf{y})$.
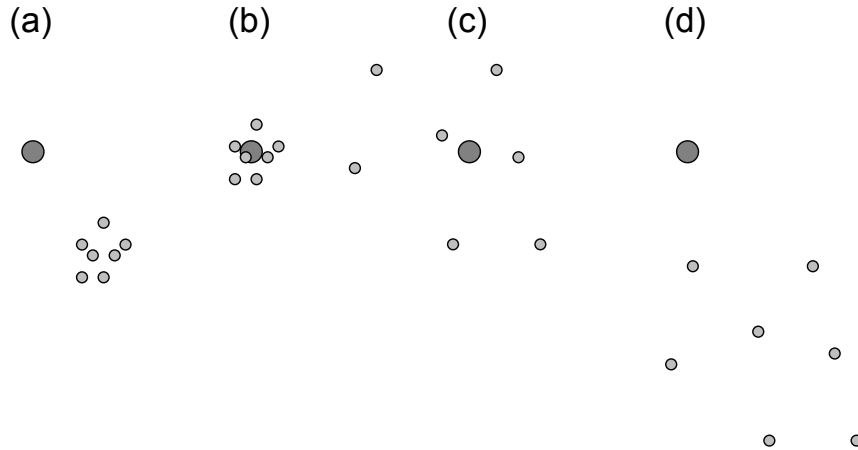
An error can be viewed as having a random component and a systematic component. The random error is different from measurement to measurement, whereas the systematic error is the same for each measurement. A measurement which has a small random uncertainty is said to have high precision, while a measurement which has a small systematic error is said to have high accuracy. This is shown conceptually in Figure 8.2.

A correction can be applied to compensate for systematic effects. It is assumed that, after correction, the expected value of the error arising from a systematic effect is zero. The correction factor is not known perfectly, so even after compensation there remains a systematic uncertainty.

The variation of random error can be described by a probability density function (pdf) such that the expected value of the random error is zero. As the random error often arises from the addition of many effects, the central limit theorem suggests that a Gaussian distribution is a good representation of this pdf. The Gaussian or Normal function with mean $\mu$ and standard deviation $\sigma$ is denoted $N(\bar{x}, \sigma)$ and defined by

$$N(\mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[ -\frac{(x - \mu)^2}{2\sigma^2} \right]. \tag{8.15}$$

Therefore, the uncertainty value commonly adopted for a scale measurand is the one-sigma standard deviation of repeated measurements of the same quantity under the

**FIGURE 8.2**
In each of the above figures the true value of the measurand is represented by the large grey dot, while repeated measurements are shown by the smaller grey dots. The measurements would be described as: a) precise but inaccurate, b) precise and accurate, c) imprecise but accurate, d) imprecise and inaccurate.

same conditions. A Gaussian pdf is usually very good at describing the uncertainty resulting from error contributions of a number of independent random processes. This can be interpreted as a result of the central limit theorem, which can be loosely summarised as *the sum of n random variates tends to a Gaussian as $n \to \infty$*. For an $m$ dimensional measurement vector the Gaussian measurement probability density function is

$$P(\mathbf{y}) = \frac{1}{\sqrt{(2\pi)^m}|\mathbf{S}_\mathrm{y}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{y} - \bar{\mathbf{y}})^\mathrm{T}\mathbf{S}_\mathrm{y}^{-1}(\mathbf{y} - \bar{\mathbf{y}})\right] \tag{8.16}$$

where $\mathbf{S}_\mathrm{y}$ is the covariance matrix of $\mathbf{y}$ and $\bar{\mathbf{y}}$ is the measurement in the absence of error.

A set of random errors $\{\epsilon_1, \epsilon_2, \ldots, \epsilon_N\}$ has a mean value of zero by definition so that the standard deviation is calculated as

$$\sigma = \sqrt{\frac{1}{N}\sum_{i=1}^{N}\epsilon_i^2} = \sqrt{\langle\epsilon^2\rangle} \tag{8.17}$$

which is also referred to as the root-mean-square or RMS value.

If $\mathbf{y}$ is an $m$ element vector with an associated error vector

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{bmatrix} \tag{8.18}$$

then the uncertainty on the measurements is represented by the measurement uncertainty covariance matrix, $\mathbf{S}_y$, defined by

$$\mathbf{S}_y = \begin{bmatrix} \text{Cov}(\epsilon_1, \epsilon_1) & \text{Cov}(\epsilon_1, \epsilon_2) & \cdots & \text{Cov}(\epsilon_1, \epsilon_m) \\ \text{Cov}(\epsilon_2 \epsilon_1) & \text{Cov}(\epsilon_2, \epsilon_2) & \cdots & \text{Cov}(\epsilon_2, \epsilon_m) \\ \vdots & \vdots & & \vdots \\ \text{Cov}(\epsilon_m, \epsilon_1) & \text{Cov}(\epsilon_m, \epsilon_2) & \cdots & \text{Cov}(\epsilon_m, \epsilon_m) \end{bmatrix}. \tag{8.19}$$

As the $\langle \epsilon_i \rangle = 0$, Equation 8.19 can be written more economically as

$$\mathbf{S}_y = \begin{bmatrix} \langle \epsilon_1 \epsilon_1 \rangle & \langle \epsilon_1 \epsilon_2 \rangle & \cdots & \langle \epsilon_1 \epsilon_m \rangle \\ \langle \epsilon_2 \epsilon_1 \rangle & \langle \epsilon_2 \epsilon_2 \rangle & \cdots & \langle \epsilon_2 \epsilon_m \rangle \\ \vdots & \vdots & & \vdots \\ \langle \epsilon_m \epsilon_1 \rangle & \langle \epsilon_m \epsilon_2 \rangle & \cdots & \langle \epsilon_m \epsilon_m \rangle \end{bmatrix} = \langle \epsilon \epsilon^{\mathrm{T}} \rangle. \tag{8.20}$$

The diagonal elements of the matrix are the variances of the corresponding element of $\mathbf{y}$. The off-diagonal elements of the uncertainty covariance matrix allow for correlation between errors. In the same way the variance represents the uncertainty in an individual value the covariance matrix represents the uncertainty in a vector.

### 8.1.3 Error Propagation

A measurement may provide involve intermediate values from which the final value of the measurand is estimated. If the covariance matrix of the intermediate values is known a method is needed to propagate the uncertainty into the dimensions of the measurand. For example suppose we want to determine the uncertainty in a quantity $y$ that is a function of two components $x_1, x_2$. To estimate the uncertainty in $y$ the mean value for $y$ is assumed to be given by the transform of the mean values of $x_1$ and $x_2$, i.e.

$$\bar{y} = f(\bar{x}_1, \bar{x}_2) \tag{8.21}$$

The spread of $y_i$ values is found from a Taylor expansion about the mean, i.e.

$$y_i - \bar{y} \approx (x_{1,i} - \bar{x}_1)\left(\frac{\partial \bar{y}}{\partial \bar{x}_1}\right) + (x_{2,i} - \bar{x}_2)\left(\frac{\partial \bar{y}}{\partial \bar{x}_2}\right) \tag{8.22}$$

The variance of $y_i$ is then given by

$$\begin{aligned} \sigma_y^2 &= \frac{1}{m} \sum_{i=1}^{m} (y_i - \bar{y})^2 = \frac{1}{m} \sum_{i=1}^{m} \left[ (x_{1,i} - \bar{x}_1)\left(\frac{\partial \bar{y}}{\partial \bar{x}_1}\right) + (x_{2,i} - \bar{x}_2)\left(\frac{\partial \bar{y}}{\partial \bar{x}_2}\right) \right]^2 \\ &= \frac{1}{m} \sum_{i=1}^{m} \left[ (x_{1,i} - \bar{x}_1)^2 \left(\frac{\partial \bar{y}}{\partial \bar{x}_1}\right)^2 + (x_{2,i} - \bar{x}_2)^2 \left(\frac{\partial \bar{y}}{\partial \bar{x}_2}\right)^2 + 2(x_{1,i} - \bar{x}_1)(x_{2,i} - \bar{x}_2)\left(\frac{\partial \bar{y}}{\partial \bar{x}_1}\right)\left(\frac{\partial \bar{y}}{\partial \bar{x}_2}\right) \right] \\ &= \sigma_{x_1}^2 \left(\frac{\partial \bar{y}}{\partial \bar{x}_1}\right)^2 + \sigma_{x_2}^2 \left(\frac{\partial \bar{y}}{\partial \bar{x}_2}\right)^2 + 2\text{Cov}(x_1, x_2)\frac{\partial \bar{y}}{\partial \bar{x}_1}\frac{\partial \bar{y}}{\partial \bar{x}_2} \end{aligned} \tag{8.23}$$

In the last line of Equation 8.23 the first two terms can be considered to be the averages of the squares of the deviations in $y$ produced by the uncertainty in $x_1$ and $x_2$ respectively. The third term accounts for and correlation between deviations in $x_1$ and $x_2$. Equation 8.23 can be generalised to a function of $N$ variables and expressed in matrix form as

$$\sigma_y^2 = \mathbf{k}\mathbf{S}_x\mathbf{k}^{\mathrm{T}} \tag{8.24}$$

where $\mathbf{S}_x = \langle \epsilon\epsilon^{\mathrm{T}} \rangle$ is the error uncertainty matrix and $\mathbf{k}$ is the vector defined by

$$k_i = \frac{\partial y}{\partial x_i} \quad \text{or} \quad \mathbf{k} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \cdots & \frac{\partial y}{\partial x_N} \end{bmatrix}. \tag{8.25}$$

For two independent parameters the covariance terms are zero so that Equation 8.23 has the simple form

$$\sigma_y^2 = \sigma_{x_1}^2 \left( \frac{\partial y}{\partial x_1} \right)^2 + \sigma_{x_2}^2 \left( \frac{\partial y}{\partial x_2} \right)^2 \qquad \text{[for no covariance]} \tag{8.26}$$

The propagation of error for a multivariate vector function cab be derived through an analogous argument. First consider $\mathbf{f} = \{f_1, f_2, \ldots, f_M\}$ a set of $M$ functions of $\mathbf{x}$ which are expanded in a Taylor series about the mean i.e.

$$f_i \approx f_i(\mu_1, \mu_2, \ldots, \mu_N) + \left( \frac{\partial f_i}{\partial x_1} \right)(x_1 - \mu_1) + \left( \frac{\partial f_i}{\partial x_2} \right)(x_2 - \mu_2) + \ldots + \left( \frac{\partial f_i}{\partial x_N} \right)(x_N - \mu_N). \tag{8.27}$$

The covariance between $f_k$ and $f_l$ is

$$\begin{aligned}
\mathrm{Cov}(f_k, f_l) &= \langle f_k f_l \rangle - \langle f_k \rangle \langle f_l \rangle \\
&= \left( \frac{\partial f_k}{\partial x_1} \right)(x_1 - \mu_1)\left( \frac{\partial f_l}{\partial x_1} \right)(x_1 - \mu_1) + \left( \frac{\partial f_k}{\partial x_1} \right)(x_1 - \mu_1)\left( \frac{\partial f_l}{\partial x_2} \right)(x_2 - \mu_2) + \ldots \\
&\qquad + \left( \frac{\partial f_k}{\partial x_1} \right)(x_1 - \mu_1)\left( \frac{\partial f_l}{\partial x_N} \right)(x_N - \mu_N) \\
&+ \left( \frac{\partial f_k}{\partial x_2} \right)(x_2 - \mu_2)\left( \frac{\partial f_l}{\partial x_1} \right)(x_1 - \mu_1) + \left( \frac{\partial f_k}{\partial x_2} \right)(x_2 - \mu_2)\left( \frac{\partial f_l}{\partial x_2} \right)(x_2 - \mu_2) + \ldots \\
&\qquad + \left( \frac{\partial f_k}{\partial x_2} \right)(x_2 - \mu_2)\left( \frac{\partial f_l}{\partial x_N} \right)(x_N - \mu_N) \\
&\quad \vdots \\
&+ \left( \frac{\partial f_k}{\partial x_N} \right)(x_N - \mu_N)\left( \frac{\partial f_l}{\partial x_1} \right)(x_1 - \mu_1) + \left( \frac{\partial f_k}{\partial x_N} \right)(x_N - \mu_N)\left( \frac{\partial f_l}{\partial x_2} \right)(x_2 - \mu_2) + \ldots \\
&\qquad + \left( \frac{\partial f_k}{\partial x_N} \right)(x_N - \mu_N)\left( \frac{\partial f_l}{\partial x_N} \right)(x_N - \mu_N) \rangle \\
&= \sum_{i=1}^{N} \sum_{j=1}^{N} \left( \frac{\partial f_k}{\partial x_i} \right)\left( \frac{\partial f_l}{\partial x_j} \right)\mathrm{Cov}(x_i, x_j).
\end{aligned}$$

If the uncertainty covariance of $\mathbf{x}$ is $\mathbf{S}_x$ then this summation is more elegantly expressed in terms of matrices

$$\mathbf{S}_y = \mathbf{K}\mathbf{S}_x\mathbf{K}^T \tag{8.28}$$

where $\mathbf{K}$ is the Jacobian of the set of $\mathbf{f}$ with respect to $\mathbf{x}$ i.e.

$$K_{ij} = \frac{\partial f_i}{\partial x_j} \quad \text{or} \quad \mathbf{K} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_2(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}. \tag{8.29}$$

The dimensions of $\mathbf{K}$ are $n \times m$ (i.e. the size of $\mathbf{x}$ times the number of functions).

## Example

Consider a series of independent measurements $\mathbf{y}$ at points $\mathbf{x}$ shown in Figure. If linear interpolation is used to estimate the value at point $c$ between $x_a$ and $x_b$ the value and the uncertainty could be estimated by standard formulae e.g.

$$y_c = y_a \frac{(x_b - x_c)}{(x_b - x_a)} + y_b \frac{(x_c - x_a)}{(x_b - x_a)} \quad \delta_{y_c} = \sqrt{\delta_{y_a}^2 \left( \frac{(x_b - x_c)}{(x_b - x_a)} \right)^2 + \delta_{y_b}^2 \left( \frac{x_c - x_a}{x_b - x_a} \right)^2} \tag{8.30}$$

Interpolation is a weighted average so the interpolate has a reduced the uncertainty, e.g. if $\delta_y = \delta_{y_a} = \delta_{y_b}$ and if $c$ lies half-way between $a$ and $b$ then

$$y_c = \frac{y_a}{2} + \frac{y_b}{2} \qquad \delta_{y_c} = \frac{\delta_y}{\sqrt{2}} \tag{8.31}$$

So then why can't the process be repeated and the error reduced by $\sqrt{2}$ each time? The point here is that successive interpolates are correlated as they have been derived from a common point. If you repeat the interpolation you need to use the full covariance matrix.

To see what the covariance matrix looks like consider the uncertainty on the $\mathbf{y}$ to be $\delta_y$ and the estimation of the mid-point values would be achieved using a Jacobian $\mathbf{K}$ defined by

$$\mathbf{K} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 & 0 & \cdots & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 & \cdots & 0 \\ \vdots & & & & & \\ 0 & 0 & \cdots & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix} \tag{8.32}$$

and the covariance matrix on the interpolated values is $\mathbf{K}^T \delta_y \mathbf{K}$.

It is straightforward to rearrange Equation 8.28 show that the transform of $\mathbf{S}_y$ to $\mathbf{S}_x$ is

$$\mathbf{S}_x^{-1} = \mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K}. \tag{8.33}$$

## 8.2 Retrieval

A retrieval problem has two distinct components: the direct (or forward) problem, which is the calculation of the measurements from the target quantity, and the inverse problem, which is the estimation of the target quantity from a set of measurements with associated errors. Much of the work involved in atmospheric remote sensing is contained in constructing *the forward model*. This is the mathematical description of the atmosphere and of the measuring instrument. The forward model describes the values that are measured as a function of the atmospheric state. The general approach involves establishing a set of forward model equations that are then solved, possibly using additional external constraints. Generally, the forward model equations are not directly invertible and so mathematical inversion techniques are employed.

In its most simple form the relationship between a measurement, $y$, and the target quantity, $x$, is represented by

$$y = f(x, \mathbf{b}) + \epsilon \tag{8.34}$$

where $f(x)$ is called the *forward function* and embodies the physics of the measurement. The measurement error is denoted by $\epsilon$ while $\mathbf{b}$ represents parameters that are needed by the forward model but are already well known so are not retrieved. A simple example would be the estimation of the temperature of a body after measurement of its emitted spectral radiance. In this case $f$ would denote the Planck function, $x$ the temperature, $y$ the measured radiance, and $\mathbf{b}$ the known wavelength, the Planck function constants and the spectral emissivity.
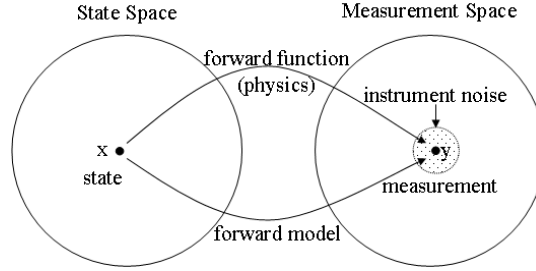
### 8.2.1 Measurement and State Space

In any problem, the finite number of measurements means that only a finite number of atmospheric properties can be determined. This statement is formalised by considering $m$ measurements denoted by $y_1, y_2, \ldots y_m$. These values are assembled into a vector, $\mathbf{y} = \{y_1, y_2, \ldots y_m\}$, which is called *the measurement vector*. The m-dimensional space of measurement vectors is called *the measurement space*. The quantities to be retrieved, $x_1, x_2, \ldots x_n$, can be assembled into an $n$ element vector, $\mathbf{x} = \{x_1, x_2, \ldots x_n\}$, which is called *the state vector*. In general, the state vector is a discrete representation of a continuous function, e.g. atmospheric temperature at a number of levels rather than a continuous profile. The $n$-dimensional space of state vectors is called *the state space*.

The multidimensional form of Equation 8.34 is

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) + \epsilon. \tag{8.35}$$

As the forward function is exact, the error term arises purely from measurement noise. The forward model, $\mathbf{F}(\mathbf{x}, \mathbf{b})$, is the computational representation of the forward function, where $\mathbf{b}$ represents a vector composed of the ancillary known parameters that are used by the forward model but are not retrieved. For simplicity, the forward

**FIGURE 8.3**

A simplified view of the relationship between state space and measurement space.

function's dependence on the ancillary parameters is not usually explicitly shown. The general representation of a measurement is therefore

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) + \boldsymbol{\epsilon} \qquad (8.36)$$

where $\boldsymbol{\epsilon}$ is not the same as in Equation 8.35 as it must include contributions from the approximation of the forward function by the forward model. A simplified representation of the terms introduced is shown in Figure 8.3. It should be noted that the mapping from state space to measurement space is not one-to-one. The subspace of state space that maps onto a point in measurement space is called null space. Null space represents a volume in state space that cannot be measured because it maps to the same point in measurement space.

The approximation of the forward function by the forward model introduces further error so that the covariance, $\mathbf{S}_{\epsilon}$, defined by

$$\mathbf{S}_{\epsilon} = \langle (\mathbf{y} - \mathbf{F}(\mathbf{x}))(\mathbf{y} - \mathbf{F}(\mathbf{x}))^{\mathrm{T}} \rangle \qquad (8.37)$$

describes the total uncertainty between $\mathbf{y}$ and $\mathbf{F}(\mathbf{x})$.

### 8.2.2   The Linear Forward Model

If the measurements are linearly related to the state (i.e. the target quantities) then Equation 8.36 can be generalised and written

$$\mathbf{y} = \mathbf{K}\mathbf{x} + \boldsymbol{\epsilon} \qquad (8.38)$$

where $\mathbf{K}$ is the Jacobian also called weighting function matrix in this context.

The reciprocal of Equation 8.38 relates the estimate of the state $\hat{\mathbf{x}}$ to the measurement vector, i.e.

$$\hat{\mathbf{x}} = \mathbf{G}\mathbf{y} \qquad (8.39)$$

where $\mathbf{G}$ is called the gain matrix. This matrix can be used to map the total error to the uncertainty in the state, i.e.

$$\mathbf{S}_x = \mathbf{G}\mathbf{S}_{\epsilon}\mathbf{G}^{\mathrm{T}}. \qquad (8.40)$$

An obvious approach (and potential mistake) is to assume that $\mathbf{G} = \mathbf{K}^{-1}$. To show how this approach can go wrong, consider two observations of a system described by two state values where the weighting function is

$$\mathbf{K} = \begin{bmatrix} 1 & 1 \\ 2 & 2.01 \end{bmatrix}. \tag{8.41}$$

If the state has the values [1, 1] then the measurement vector is

$$\mathbf{y} = \mathbf{Kx} = \begin{bmatrix} 1 & 1 \\ 2 & 2.01 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 4.01 \end{bmatrix}. \tag{8.42}$$

The inverse of $\mathbf{K}$ [$= \mathbf{G}$ in this case] is

$$\mathbf{K}^{-1} = \begin{bmatrix} 201 & -100 \\ -200 & 100 \end{bmatrix}. \tag{8.43}$$

This matrix returns the correct state, [1, 1], if the measurement vector is [2,4.01]. However, if one of the measurements in the measurement vector is slightly wrong, i.e. $\mathbf{y} = [2, 4]$, then the solution would change to

$$\hat{\mathbf{x}} = \mathbf{Gy} = \begin{bmatrix} 201 & -100 \\ -200 & 100 \end{bmatrix} \begin{bmatrix} 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}. \tag{8.44}$$

There is a 100 % difference in each solution for a 0.25 % change in a single measurement. More generally, if the measurements are independent with a total error of 10 % of the signal, i.e. $[0.2, 0.4]$, then the error in the solution is

$$\mathbf{S}_{\hat{\mathbf{x}}} = \mathbf{GS}_\epsilon \mathbf{G}^{\mathrm{T}} = \begin{bmatrix} 201 & -100 \\ -200 & 100 \end{bmatrix} \begin{bmatrix} 0.04 & 0 \\ 0 & 0.16 \end{bmatrix} \begin{bmatrix} 201 & -200 \\ -100 & 100 \end{bmatrix} \approx \begin{bmatrix} 3216 & -3208 \\ -3208 & 3200 \end{bmatrix}. \tag{8.45}$$

If the uncertainty on the state is estimated by the square root of the diagonal of $\mathbf{S}_x$ then the solutions for $x_1$ and $x_2$ are $2 \pm 56.7$ and $4 \pm 56.6$ respectively. The magnification of error, and so lack of precision in the result, provides motivation to find an alternative retrieval approach.

### 8.2.3 Types of Inverse Problem

Assuming $\mathbf{K}$ is of full rank (the rank of a matrix corresponds to the number of linearly independent rows or columns of the matrix) there are three types of problem to consider:

- underconstrained $m < n$

- well-posed $m = n$

- overconstrained $m > n$

For underconstrained problems there are two approaches:

- reduce the number of unknowns from infinity to, at most, *m*. This can be achieved by parameterising the state in some way, e.g. by representing a temperature profile with an *m*- term Fourier expansion.

- add $n - m$ constraints.

For overconstrained or well-posed problems a standard inverse method is that of the least-squares solution.

The determination of the gain matrix and hence the solution of the inverse problem can be approach in two ways. These are:

- finding an exact solution where the forward modelled solution ($\mathbf{K\hat{x}}$) matches the measurements exactly, or

- finding a non-exact solution where the forward modelled solution fits 'within experimental error'.

These two approaches are explored in the following section in the application of the least-squares solution to the inverse problem.

## 8.3   Least-Squares Solution

A common approach to finding a solution is to adjust the state in order to minimise the sum-square difference, $\chi^2$, between the measurements and the predicted measurements $\mathbf{F}(\mathbf{x})$. The value of $\chi^2$ is found from

$$\chi^2 = \sum_{i=1}^{m} \left( y_i - \sum_{j=1}^{n} K_{ij} x_j \right)^2 = (\mathbf{y} - \mathbf{Kx})^{\mathrm{T}} (\mathbf{y} - \mathbf{Kx}). \tag{8.46}$$

To find the minimum value of $\chi^2$ we find where the derivative is zero, i.e.

$$\frac{\partial}{\partial x_k} \sum_{i=1}^{m} \left( y_i - \sum_{j=1}^{n} K_{ij} x_j \right)^2 = \sum_{i=1}^{m} 2K_{ik} \left( y_i - \sum_{j=1}^{n} K_{ij} x_j \right)$$

$$= 2 \left( \sum_{i=1}^{m} K_{ik} y_i - \sum_{i=1}^{m} K_{ik} \sum_{j=1}^{n} K_{ij} x_j \right) = 0 \tag{8.47}$$

which can be expressed as

$$\sum_{i=1}^{m} \sum_{j=1}^{n} K_{ki}^{\mathrm{T}} K_{ij} x_j = \sum_{i=1}^{m} K_{ki}^{\mathrm{T}} y_i. \tag{8.48}$$

The least-squares solution is then

$$x_j = \frac{\sum_{i=1}^{m} K_{ki}^{\mathrm{T}} y_i}{\sum_{i=1}^{m} K_{ki}^{\mathrm{T}} K_{ij}} \text{ for } j = 1 \text{ to } n. \tag{8.49}$$

Performing the derivation in matrix notation is quicker:

$$\frac{\partial}{\partial x} \left[ (\mathbf{y} - \mathbf{Kx})^{\mathrm{T}} (\mathbf{y} - \mathbf{Kx}) \right] = -2\mathbf{K}^{\mathrm{T}} (\mathbf{y} - \mathbf{Kx}) = -2\mathbf{K}^{\mathrm{T}}\mathbf{y} + 2\mathbf{K}^{\mathrm{T}}\mathbf{Kx} = 0$$

$$\Rightarrow \mathbf{K}^{\mathrm{T}}\mathbf{Kx} = \mathbf{K}^{\mathrm{T}}\mathbf{y}$$

$$\Rightarrow \mathbf{x} = (\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{y}. \tag{8.50}$$

From this expression it is clear that the gain matrix is

$$\mathbf{G} = (\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}. \tag{8.51}$$

Substituting this into the error propagation expression (Equation 8.40) gives the co-variance error matrix for the least-squares solution as

$$\mathbf{S}_x = (\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon \left[ (\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}} \right]^{\mathrm{T}} = (\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon\mathbf{K}(\mathbf{K}^{\mathrm{T}}\mathbf{K})^{-1} \tag{8.52}$$

which makes use of the standard matrix expression $(\mathbf{AB})^{\mathrm{T}} = \mathbf{B}^{\mathrm{T}}\mathbf{A}^{\mathrm{T}}$ for matrices $\mathbf{A}$ and $\mathbf{B}$. Note that as $\mathbf{K}^{\mathrm{T}}\mathbf{K}$ is symmetric then so also is its inverse (and a symmetric matrix is unchanged by the transpose operation).

## 8.3.1  Least-Squares Solution with Errors

If the uncertainty in the measurements and other errors are included in the expression for $\chi^2$ then

$$\chi^2 = (\mathbf{y} - \mathbf{Kx})^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{Kx}). \tag{8.53}$$

Finding the minimum as before

$$\frac{\partial}{\partial x} \left[ (\mathbf{y} - \mathbf{Kx})^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{Kx}) \right] = -2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{Kx}) = -2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{y} + 2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{Kx} = 0$$

$$\Rightarrow \mathbf{x} = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{y} \tag{8.54}$$

with an associated gain matrix

$$\mathbf{G} = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1} \tag{8.55}$$

and expression for the error covariance of the state, i.e.

$$\mathbf{S}_x = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{S}_\epsilon \left[ (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1} \right]^{\mathrm{T}} = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K})^{-1}. \tag{8.56}$$

which is nothing more than a re-expression of Equation 8.33 which defines the transform of the covariance of $y$ to the covariance of $x$.
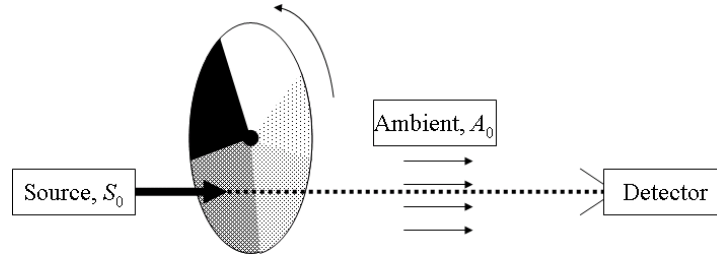
**FIGURE 8.4**
An example radiative retrieval problem.

The result can be constrained by the use of a priori information, $\mathbf{x}_a$, which is the pre-measurement knowledge of the state with an associated uncertainty, $\mathbf{S}_a$. The introduction of the constraint changes the cost function to

$$\chi^2 = (\mathbf{y} - \mathbf{Kx})^T \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{Kx}) + (\mathbf{x} - \mathbf{x}_a)^T \mathbf{S}_a^{-1} (\mathbf{x} - \mathbf{x}_a) \qquad (8.57)$$

where the first term is the least-squares cost while the second term is the cost associated with the constraint.

A potential problem with the method of least-squares (and other approaches) is that $\mathbf{K}$ may be singular. Remember, matrix $\mathbf{K}$ is singular if $\mathbf{K}^{-1}$ does not exist. Even if $\mathbf{K}$ is non-singular, it may be ill-conditioned. A matrix is ill-conditioned if it is invertible but can become singular (non-invertible) if some of its entries are changed slightly. This was the case for $\mathbf{K}$ in the example given at the end of Section 8.2.2.

### 8.3.2   An Example Retrieval Problem

To illustrate the set-up of retrieval problems consider the situation shown in Figure 8.4. A detector views a source of magnitude $S_0$ in the presence of an ambient signal, $A_0$. The source is modulated by a disk that is divided into five sections of differing opacity. The detector is timed to give five measurements when the source radiates through each of the opaque regions. The transmissions $T$ of the opaque regions of the disk are known to be 0, 0.25, 0.5, 0.75 and 1 respectively. The problem is to determine the source strength and the ambient signal from the five measurements.

In this situation the measurement vector is

$$\mathbf{y} = \begin{bmatrix} A_0 \\ 0.25S_0 + A_0 \\ 0.5S_0 + A_0 \\ 0.75S_0 + A_0 \\ S_0 + A_0 \end{bmatrix} = \begin{bmatrix} 2.1 \\ 30.5 \\ 42.7 \\ 79.2 \\ 90.5 \end{bmatrix} \qquad (8.58)$$

where the 'true' values of $S_0$ and $A_0$ have been taken as 100 and 2, and the measurements have been simulated assuming a measurement noise of 10 % and that the forward model is exact. These values are in some arbitrary units that will be left

unspecified. The error has been taken to be uncorrelated so that the total error co-variance matrix is

$$\mathbf{S}_\epsilon = \mathbf{S}_y = \begin{bmatrix} 0.04 & 0 & 0 & 0 & 0 \\ 0 & 7.29 & 0 & 0 & 0 \\ 0 & 0 & 27.04 & 0 & 0 \\ 0 & 0 & 0 & 59.29 & 0 \\ 0 & 0 & 0 & 0 & 104.04 \end{bmatrix}. \tag{8.59}$$

The weighting function (which is generally found numerically) has an analytic form for this example. It is

$$\mathbf{K} = \begin{bmatrix} 0 & 1 \\ 0.25 & 1 \\ 0.5 & 1 \\ 0.75 & 1 \\ 1 & 1 \end{bmatrix}. \tag{8.60}$$

The gain matrix is evaluated as

$$\mathbf{G} = (\mathbf{K}^T\mathbf{K})^{-1}\mathbf{K}^T = \begin{bmatrix} -0.75 & -0.44 & .015 & 0.40 & 0.78 \\ 0.57 & 0.42 & 0.19 & 0.0046 & -0.18 \end{bmatrix} \tag{8.61}$$

and the solution is

$$\mathbf{x} = \mathbf{G}\mathbf{y} = \begin{bmatrix} 87.9 \\ 5.9 \end{bmatrix}. \tag{8.62}$$

Propagating this uncertainty into state space (Equation 8.52) gives the uncertainty on our estimates as

$$\mathbf{S}_x = \begin{bmatrix} 74.6 & -16.1 \\ -16.1 & 5.78 \end{bmatrix}. \tag{8.63}$$

If the least-squares solution includes the uncertainty in the fits then the gain matrix is

$$\mathbf{G} = (\mathbf{K}^T\mathbf{K})^{-1}\mathbf{K}^T = \begin{bmatrix} -2.01 & -0.80 & 0.55 & 0.38 & 0.29 \\ 1.00 & 0.0033 & -0.00001 & -0.0003 & -0.0004 \end{bmatrix} \tag{8.64}$$

and the solution (from Equation 8.54) is

$$\mathbf{x} = \mathbf{G}\mathbf{y} = \begin{bmatrix} 99.0 \\ 2.1 \end{bmatrix} \tag{8.65}$$

with the associated error covariance

$$\mathbf{S}_x = (\mathbf{K}^T\mathbf{S}_y^{-1}\mathbf{K})^{-1} = \begin{bmatrix} 29.7 & -0.08 \\ -0.08 & 0.040 \end{bmatrix}. \tag{8.66}$$

There is a clear improvement in fit and a reduction in uncertainty of the state in this solution.

## 8.4   Optimal Estimation Solution

There is usually an infinite set of all possible solutions consistent with the measurement vector and its associated errors. In optimal estimation, prior knowledge of the state is used to constrain the possible solutions. The prior knowledge may come from earlier measurements, model output or even an educated guess. The use of a priori information is formalised through Bayes' theorem.

### 8.4.1   Bayes' Theorem

Bayes' theorem can be derived by considering the probability of event $a$, $P(a)$, and event $b$, $P(b)$. The probability of event $a$ occurring given that event $b$ occurs is the conditional probability $P(a|b)$ and is defined

$$P(a|b) = \frac{P(a \cap b)}{P(b)} \tag{8.67}$$

where $P(a \cap b)$ represents the probability of both $a$ and $b$ occurring. The inverse conditional probability $P(b|a)$ can similarly be represented

$$P(b|a) = \frac{P(a \cap b)}{P(a)}. \tag{8.68}$$

Rearranging these equations gives

$$P(a|b)P(b) = P(a \cap b) = P(b|a)P(a)$$
$$\Rightarrow P(a|b) = \frac{P(b|a)P(a)}{P(b)}. \tag{8.69}$$

As an example of the application of Bayes' theorem consider two bags. In one bag there are 3 blue marbles and 6 red, while in the second bag there is one blue marble and 4 red. A die is tossed and if the numbers 1 or 2 are thrown a marble is selected at random from the first bag, otherwise a marble is selected at random from the second bag. The probability of getting a blue marble is $\frac{1}{3} \times \frac{3}{9} + \frac{2}{3} \times \frac{1}{5} = \frac{11}{45}$. Bayes' theorem can be used to answer the question "*If a blue marble is drawn then what is the probability that the selection was from bag 1?*" as follows

$$P(\text{Bag 1}|\text{Blue}) = \frac{P(\text{Blue}|\text{Bag 1})P(\text{Bag 1})}{P(\text{Blue})} = \frac{\frac{3}{9} \times \frac{1}{3}}{\frac{11}{45}} = \frac{5}{11}. \tag{8.70}$$

When applied to a retrieval problem, Bayes' theorem defines the probability density function of the state given the measurement, i.e.

$$P(\mathbf{x}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{x})P(\mathbf{x})}{P(\mathbf{y})}. \tag{8.71}$$

where the value of $P(\mathbf{y})d\mathbf{y}$ is the probability that the measurement was in the multidimensional interval $(\mathbf{y}, \mathbf{y} + d\mathbf{y})$. Similarly, the probability density function, $P(\mathbf{x})$, has values such that $P(\mathbf{x})d\mathbf{x}$ is the probability that the true value of the state lies in the interval $(\mathbf{x}, \mathbf{x} + d\mathbf{x})$.

However, describing the solution of a problem by a pdf is not usually very helpful. A criterion is needed to select a solution from the constrained set of possible solutions. The criteria that could be adopted include

- the most probable solution, $\mathbf{x}_{\text{max}}$, given by

$$\left.\frac{dP(\mathbf{x})}{d\mathbf{x}}\right|_{\mathbf{x}_{\text{max}}} = 0 \tag{8.72}$$

- the expected solution, $\langle\mathbf{x}\rangle$, given by

$$\langle\mathbf{x}\rangle = \int P(\mathbf{x})\mathbf{x}\,d\mathbf{x} \tag{8.73}$$

where the integral represents multidimensional integration over all of state space.

## 8.4.2 The Linear Retrieval Problem

If Gaussian statistics are used to describe the prior knowledge of the state then $P(\mathbf{x})$ is given by

$$P(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n}|\mathbf{S}_{\text{a}}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{x}_{\text{a}})^{\text{T}}\mathbf{S}_{\text{a}}^{-1}(\mathbf{x} - \mathbf{x}_{\text{a}})\right] \tag{8.74}$$

where $\mathbf{x}_{\text{a}}$ is the a priori value of $\mathbf{x}$ and $\mathbf{S}_{\text{a}}$ is the associated covariance matrix.
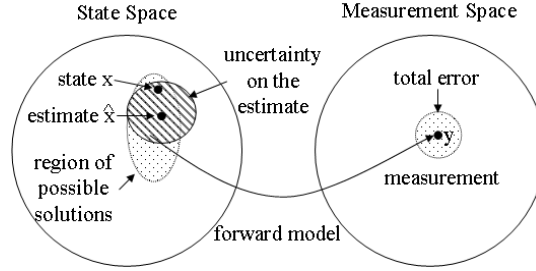
The act of measurement maps the state into measurement space using the forward function. This transform is represented by the forward model, $\mathbf{F}(\mathbf{x})$, so the conditional probability of measuring $\mathbf{y}$ given state $\mathbf{x}$ is

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^m}|\mathbf{S}_\epsilon|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{F}(\mathbf{x}))^{\text{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{F}(\mathbf{x}))\right] \tag{8.75}$$

where $\mathbf{S}_\epsilon$ is used instead of $\mathbf{S}_{\text{y}}$ to include error introduced by approximating the forward function by the forward model.

In the inverse problem, it is recognised that a given measurement $\mathbf{y}$ could be the result of a mapping from within a region of state space rather than from a single point. This idea is shown graphically in Figure 8.5. Bayes' theorem is used to identify the region of state space in which the solution lies. The region is identified using the conditional probability of $\mathbf{x}$ given $\mathbf{y}$, i.e.

$$P(\mathbf{x}|\mathbf{y}) = \left\{\frac{1}{\sqrt{(2\pi)^m}|\mathbf{S}_\epsilon|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{F}(\mathbf{x}))^{\text{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{F}(\mathbf{x}))\right]\right.$$

$$\left.\times \frac{1}{\sqrt{(2\pi)^n}|\mathbf{S}_{\text{a}}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{x}_{\text{a}})^{\text{T}}\mathbf{S}_{\text{a}}^{-1}(\mathbf{x} - \mathbf{x}_{\text{a}})\right]\right\}/P(\mathbf{y}). \tag{8.76}$$

**FIGURE 8.5**
On applying Bayes' theorem to the retrieval problem a region in state space is identified which has a finite probability of giving rise to the measurement. A value is chosen from this region as the solution.

This expression can be simplified by recognising that $P(\mathbf{y})$ is a constant for a given measurement as are the determinants of $\mathbf{S}_\epsilon$ and $\mathbf{S}_a$. In the linear case $\mathbf{F}(\mathbf{x}) = \mathbf{K}\mathbf{x}$ giving

$$P(\mathbf{x}|\mathbf{y}) = c_1 \exp\left[-\frac{1}{2}(\mathbf{y} - \mathbf{K}\mathbf{x}))^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{K}\mathbf{x}))\right]\exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{x}_a)^{\mathrm{T}}\mathbf{S}_a^{-1}(\mathbf{x} - \mathbf{x}_a)\right] \quad (8.77)$$

where

$$c_1 = \frac{1}{\sqrt{(2\pi)^{m+n}}|\mathbf{S}_\epsilon|^{1/2}|\mathbf{S}_a|^{1/2}P(\mathbf{y})}. \quad (8.78)$$

Taking the log of both sides of Equation 8.77 and slightly rearranging gives

$$-2\ln P(\mathbf{x}|\mathbf{y}) = (\mathbf{y} - \mathbf{K}\mathbf{x})^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{K}\mathbf{x}) + (\mathbf{x} - \mathbf{x}_a)^{\mathrm{T}}\mathbf{S}_a^{-1}(\mathbf{x} - \mathbf{x}_a) + c_2 \quad (8.79)$$

where

$$c_2 = -2\ln c_1. \quad (8.80)$$

The pdf of solutions in state space can be reduced to the most probable solution or the expected solution, as follows.

### 8.4.2.1 The Most Probable Solution

The most probable solution $\hat{\mathbf{x}}$ is found by minimizing the right-hand side of Equation 8.79, i.e.

$$\frac{\partial}{\partial \mathbf{x}}\left[(\mathbf{y} - \mathbf{K}\mathbf{x})^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{K}\mathbf{x}) + (\mathbf{x} - \mathbf{x}_a)^{\mathrm{T}}\mathbf{S}_a^{-1}(\mathbf{x} - \mathbf{x}_a)\right] = 0$$

$$-2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{y} - \mathbf{K}\hat{\mathbf{x}}) + 2\mathbf{S}_a(\hat{\mathbf{x}} - \mathbf{x}_a) = 0$$

$$-2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{y} + 2\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K}\hat{\mathbf{x}} + 2\mathbf{S}_a\hat{\mathbf{x}} - 2\mathbf{S}_a\mathbf{x}_a = 0.$$

Hence

$$(\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} + \mathbf{S}_a)\hat{\mathbf{x}} = \mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{y} + \mathbf{S}_a\mathbf{x}_a$$

$$\Rightarrow \hat{\mathbf{x}} = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} + \mathbf{S}_a)^{-1}(\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{y} + \mathbf{S}_a\mathbf{x}_a). \quad (8.81)$$

This is exactly the same minimization that would occur in finding the least-squared solution (with error) including a constraint as defined in Equation 8.57.

### 8.4.2.2 The Expected Solution

The expected solution is found by recognising that the product of two Gaussians is itself a Gaussian. So taking a Gaussian defined by vector mean, $\hat{\mathbf{x}}$, and covariance matrix $\hat{\mathbf{S}}$,

$$-2 \ln P(\mathbf{x}|\mathbf{y}) = (\mathbf{x} - \hat{\mathbf{x}})^{\mathrm{T}} \hat{\mathbf{S}}^{-1} (\mathbf{x} - \hat{\mathbf{x}}) + c_3 \tag{8.82}$$

and equating the quadratic terms in $\mathbf{x}$ with Equation 8.79 gives

$$\mathbf{x}^{\mathrm{T}} \mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} \mathbf{x} + \mathbf{x}^{\mathrm{T}} \mathbf{S}_{\mathrm{a}}^{-1} \mathbf{x} = \mathbf{x}^{\mathrm{T}} \hat{\mathbf{S}}^{-1} \mathbf{x}$$
$$\Rightarrow \hat{\mathbf{S}}^{-1} = \mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} + \mathbf{S}_x^{-1}. \tag{8.83}$$

Equating the linear terms in $\mathbf{x}^{\mathrm{T}}$ gives

$$-(\mathbf{K}\mathbf{x})^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{y} - \mathbf{x}^{\mathrm{T}} \mathbf{S}_{\mathrm{a}}^{-1} \mathbf{x}_{\mathrm{a}} = -\mathbf{x}^{\mathrm{T}} \hat{\mathbf{S}}^{-1} \hat{\mathbf{x}}. \tag{8.84}$$

Substituting in the expression for $\hat{\mathbf{S}}^{-1}$ (Equation 8.83) gives

$$\mathbf{x}^{\mathrm{T}} \mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{y} + \mathbf{x}^{\mathrm{T}} \mathbf{S}_{\mathrm{a}}^{-1} \mathbf{x}_{\mathrm{a}} = \mathbf{x}^{\mathrm{T}} (\mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} + \mathbf{S}_x^{-1}) \hat{\mathbf{x}}. \tag{8.85}$$

Then rearranging to make $\hat{\mathbf{x}}$ the subject of the equation gives

$$\hat{\mathbf{x}} = (\mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} + \mathbf{S}_x^{-1})^{-1} (\mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{y} + \mathbf{S}_{\mathrm{a}}^{-1} \mathbf{x}_{\mathrm{a}}) \tag{8.86}$$

which is the expected solution. It can also be expressed as

$$\hat{\mathbf{x}} = \mathbf{x}_{\mathrm{a}} + (\mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} + \mathbf{S}_{\mathrm{a}}^{-1})^{-1} \mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{K} \mathbf{x}_{\mathrm{a}}). \tag{8.87}$$

### 8.4.2.3 The Optimal Estimate Solution

Comparing the expressions for the most likely solution with the expected solution shows that they are the same, and can be found by minimizing a least-squares cost function with a Gaussian constraint on the state. This is also called the optimal estimate solution. It is convenient to express the solution relative to a reference value $\mathbf{x}_0$ in state space. As the problem is linear, $\mathbf{x}_0 = \mathbf{G}\mathbf{K}\mathbf{x}_0$ and the estimate of the state, $\hat{\mathbf{x}}$, can be expressed as

$$\hat{\mathbf{x}} = \mathbf{G}\mathbf{y} = \mathbf{x}_0 + \mathbf{G}\mathbf{y} - \mathbf{x}_0 = \mathbf{x}_0 + \mathbf{G}\mathbf{y} - \mathbf{G}\mathbf{K}\mathbf{x}_0 = \mathbf{x}_0 + \mathbf{G}(\mathbf{y} - \mathbf{K}\mathbf{x}_0). \tag{8.88}$$

Comparing this to Equation 8.87 shows that the optimal estimate can be considered as the solution expressed relative to the a priori, and the expression for the gain matrix is

$$\mathbf{G} = (\mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} \mathbf{K} + \mathbf{S}_{\mathrm{a}}^{-1})^{-1} \mathbf{K}^{\mathrm{T}} \mathbf{S}_\epsilon^{-1}. \tag{8.89}$$

Using the gain matrix, the optimal estimation solution (Equation 8.87) can be re-formed as

$$\hat{\mathbf{x}} = \mathbf{x}_{\mathrm{a}} + \mathbf{G}(\mathbf{y} - \mathbf{K}\mathbf{x}_{\mathrm{a}}) = \mathbf{x}_{\mathrm{a}} + \mathbf{G}\mathbf{y} - \mathbf{G}\mathbf{K}\mathbf{x}_{\mathrm{a}} = \mathbf{G}\mathbf{y} + (\mathbf{I}_n - \mathbf{G}\mathbf{K})\mathbf{x}_{\mathrm{a}} \tag{8.90}$$

and seen to be a weighted sum of the measurements and the a priori knowledge of the state. The error covariance matrix of the optimal estimate is given by

$$\hat{\mathbf{S}}^{-1} = \mathbf{G}\mathbf{S}_\epsilon\mathbf{G}^{\mathrm{T}} = (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} + \mathbf{S}_{\mathrm{a}}^{-1})^{-1} \tag{8.91}$$

which can be interpreted as a covariance matrix generated by the product of two Gaussian distributions: one the error covariance transformed into state space and the other the a priori error covariance.

## 8.5 Diagnostics

Although the optimal estimate solution is the most likely and the most probable solution it may not be a very good solution. One of the significant advantages of the optimal estimation approach is the ability to quantify the quality of the solution. This is done using three diagnostic measures: the averaging kernel matrix, the degrees of freedom and the information content of the measurement.

### 8.5.1 Averaging Kernel Matrix

The averaging kernel matrix, $\mathbf{A}$, relates the most probable state to the true state and is defined

$$\mathbf{A} = \frac{\partial\hat{\mathbf{x}}}{\partial\mathbf{x}}. \tag{8.92}$$

It is found by replacing $\mathbf{y}$ by $\mathbf{Kx}$ in Equation 8.87 and evaluating the derivative, i.e.

$$\begin{aligned}
\mathbf{A} = \frac{\partial\hat{\mathbf{x}}}{\partial\mathbf{x}} &= \frac{\partial}{\partial\mathbf{x}}\left[\mathbf{x}_{\mathrm{a}} + (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} + \mathbf{S}_{\mathrm{a}}^{-1})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}(\mathbf{Kx} - \mathbf{Kx}_{\mathrm{a}})\right] \\
&= (\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} + \mathbf{S}_{\mathrm{a}}^{-1})^{-1}\mathbf{K}^{\mathrm{T}}\mathbf{S}_\epsilon^{-1}\mathbf{K} \\
&= \mathbf{GK}
\end{aligned} \tag{8.93}$$

where the last simplification makes use of the identity given by Equation 8.89. Combing Equations 8.88 and 8.93 gives an expression for the optimal solution in terms of the averaging kernel matrix, i.e.

$$\hat{\mathbf{x}} = \mathbf{Ax} + (\mathbf{I} - \mathbf{A})\mathbf{x}_{\mathrm{a}} \tag{8.94}$$

so that it is possible to interpret $\mathbf{A}$ as the weights of a weighted mean of the true state and an a priori estimate that give rise to $\hat{\mathbf{x}}$. The rows of $\mathbf{A}$ are the averaging kernels (smoothing functions) that map the true state into retrieval space. The width of the kernel is a measure of retrieval resolution. The area of a kernel is approximately unity for an accurate retrieval and can be interpreted as a measure of the fraction of the retrieved value that originates from the signal and not from the a priori.

### 8.5.2  Degrees of Freedom

The number of degrees of freedom, $d$, is a scalar measure of the number of independent quantities that can be measured. For example, consider choosing three numbers at random. This process has three degrees of freedom. The degrees of freedom is reduced by each independent constraint we place on the process. For example, if the three numbers are required to add up to some number $t$ then the first two numbers can be selected at random, but the third must be chosen so that it makes the total equal to $t$ — thus the degree of freedom is two.

Now consider a set of $n$ values $\{x_1, x_2 \ldots x_n\}$ drawn from a Gaussian distribution with mean 0 and standard deviation 1. By definition, the $\chi^2$ distribution is the sum of the square of $x_n$. The expected value of $\chi^2$ is $n$, the number of degrees of freedom. If a distribution has a non-zero mean $\mu$ and a non-unity standard deviation $\sigma$ then the values are transformed to $N(0, 1)$ to get $\chi^2$, i.e.

$$\chi^2 = \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{\sigma^2}. \tag{8.95}$$

For a multivariate Gaussian distribution $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ the expression for $\chi^2$ becomes

$$\chi^2 = (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \, \boldsymbol{\Sigma}^{-1} \, (\mathbf{x} - \boldsymbol{\mu}). \tag{8.96}$$

Hence at the optimal solution the value of $\chi^2$ is expected to be

$$\chi^2 = (\mathbf{y} - \mathbf{Kx})^{\mathrm{T}} \, \mathbf{S}_{\epsilon}^{-1} \, (\mathbf{y} - \mathbf{Kx}) + (\mathbf{x} - \mathbf{x}_{\mathrm{a}})^{\mathrm{T}} \mathbf{S}_{\mathrm{a}}^{-1} (\mathbf{x} - \mathbf{x}_{\mathrm{a}}) \approx m \tag{8.97}$$

where it is assumed $m = n$. More generally, the largest number of degrees of freedom possible is determined by the number of elements in the measurement vector $m$ or the state vector $n$, whichever is the smaller. The number of degrees of freedom can be divided into two parts; the degrees of freedom for signal, $d_s$, and the degrees of freedom for noise, $d_n$. Formally, these can be expressed as [*Rodgers*, 2000]

$$d_s = \mathrm{tr}(\mathbf{A}) \tag{8.98}$$

$$d_n = \mathrm{tr}[\mathbf{S}_{\epsilon}(\mathbf{KS}_{\mathrm{a}}\mathbf{K}^{\mathrm{T}} + \mathbf{S}_{\epsilon})^{-1}]. \tag{8.99}$$

It follows that $d_s + d_n = m$. If the degrees of freedom for signal at the solution is such that $d_s < n$ then it may be that there is little information in the measurement about one of the state vector elements, or that two or more of the state vector elements are highly correlated.

### 8.5.3  Information Content of a Measurement

An obvious measure of the quality of the retrieval is given by the diagonal elements of the error covariance matrix, $\hat{\mathbf{S}}$. As a first order approximation the square root of the diagonal elements gives the retrieval error on the retrieved parameter. This value can be compared to the corresponding diagonal element in the a priori covariance matrix

to see how much the uncertainty in a parameter has been reduced by a measurement. More formally, this comparison is done using the concept of information content.

The information content of a measurement can be defined qualitatively as the factor by which knowledge of a quantity is improved by making a measurement. The Shannon definition of information content [*Shannon and Weaver*, 1949] describes reduction of entropy caused by taking a measurement. They define the entropy $S(P)$ of pdf $P$ as

$$S(P) = \sum_{i=1}^{\infty} p_i \log_2 p_i \qquad (8.100)$$

where $p_i$ is the probability of the system being in state $i$. For an $n$ dimensional Gaussian distribution whose covariance is $\mathbf{S}_y$ this evaluates to [*Rodgers*, 2000]

$$S(P) = n \log_2 \sqrt{2\pi e} + \frac{1}{2} \log_2 |\mathbf{S}_\epsilon|. \qquad (8.101)$$

If $P(x)$ is the pdf of the state before measurement and $P(x|y)$ is the pdf of the state after measurement then the Shannon information content $H$ is given by

$$H = S[P(x)] - S[P(x|y)] \qquad (8.102)$$

which for the linear Gaussian case is

$$\begin{aligned} H &= \frac{1}{2} \ln |\mathbf{S}_a| - \frac{1}{2} \ln |\hat{\mathbf{S}}| \\ &= \frac{1}{2} \log_2 |\hat{\mathbf{S}}^{-1}\mathbf{S}_a| \\ &= \frac{1}{2} \log_2 |\mathbf{I} - \mathbf{A}|. \end{aligned} \qquad (8.103)$$

Therefore the measurement will approach maximum information content as the averaging kernel approaches the identity matrix. As $\mathbf{A}$ is a readily available diagnostic the information content of a retrieval can be simply derived.

## 8.6 Non-Linearity

If the forward model is linear and the a priori state is described by Gaussian statistics then the cost function (Equation 8.57) will be quadratic in the state vector and the equations to be solved will be linear. However, in most atmospheric problems the forward model is rarely linear. For a non-linear problem it is possible to linearize the

forward model (Equation 8.38) about a reference state, $\mathbf{x}_0$, using a Taylor expansion,

$$
\begin{aligned}
\mathbf{y} - F(\mathbf{x}_0) &= (\mathbf{x} - \mathbf{x}_0) \left.\frac{\partial \mathbf{F}}{\partial \mathbf{x}}\right|_{\mathbf{x}=\mathbf{x}_0} + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^2 \left.\frac{\partial^2 \mathbf{F}}{\partial \mathbf{x}^2}\right|_{\mathbf{x}=\mathbf{x}_0} + \ldots + \boldsymbol{\epsilon} \\
&\approx (\mathbf{x} - \mathbf{x}_0) \left.\frac{\partial \mathbf{F}}{\partial \mathbf{x}}\right|_{\mathbf{x}=\mathbf{x}_0} + \boldsymbol{\epsilon} \\
&= \mathbf{K}(\mathbf{x} - \mathbf{x}_0) + \boldsymbol{\epsilon}
\end{aligned}
\tag{8.104}
$$

where the Jacobian matrix is estimated at the reference state, i.e.

$$
K_{ij} = \left.\frac{\partial F_i(\mathbf{x})}{\partial x_j}\right|_{\mathbf{x}=\mathbf{x}_0}.
\tag{8.105}
$$

To obtain a solution an initial solution is assumed: *the first guess*. The forward model is then linearised about this state and a new estimate of the solution is formed. This process continues until some convergence criterion is reached.

## 8.6.1  Iterating to a Solution

The optimal estimate solution, $\hat{\mathbf{x}}$, minimises a joint cost function $\chi^2$ defined by

$$
\chi^2 = (\mathbf{y} - F(\hat{\mathbf{x}}))^{\mathrm{T}} \mathbf{S}_\epsilon^{-1} (\mathbf{y} - F(\hat{\mathbf{x}})) + (\hat{\mathbf{x}} - \mathbf{x}_a)^{\mathrm{T}} \mathbf{S}_a^{-1} (\hat{\mathbf{x}} - \mathbf{x}_a).
\tag{8.106}
$$

There are many approaches to functional minimisation (e.g. see *Press* [1992]).
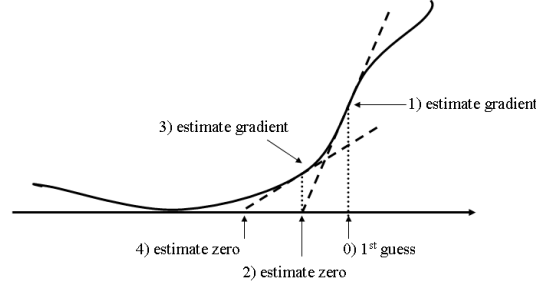
### 8.6.1.1  Steepest Descent

The method of steepest descent is an algorithm for finding the nearest local minimum of a function. A real valued function $g(\mathbf{x})$ shows the greatest change in the direction of the steepest gradient. Hence a sufficiently small step size $\gamma$ will result in a decrease in the value of the function, i.e.

$$
g(\boldsymbol{\xi} - \frac{1}{\gamma}\nabla g(\boldsymbol{\xi})) < g(\boldsymbol{\xi})
\tag{8.107}
$$

To find the optimal solution the method of steepest descent takes the form of iterating

$$
\mathbf{x}_{i+1} = \boldsymbol{\xi} - \frac{1}{\gamma}\nabla\chi^2(\boldsymbol{\xi})
\tag{8.108}
$$

starting from an initial guess, $\mathbf{x}_0$. At each step the value of $\chi^2(x_{i+1})$ is compared to the value of $\chi^2(x_i)$. If the new value of the function is smaller than the old the step then it is accepted and the value of $\gamma$ decreased. If the minimum has been 'overstepped' then the old value of $x$ is retained and $\gamma$ is increased. The initial value of $\gamma$ is usually chosen to give a reasonably small change in $x$ under the assumption that we are starting reasonably close to the solution.

**FIGURE 8.6**

Successive approximation steps to find $f(x) = 0$ using Newtonian iteration.

### 8.6.1.2  Newton's Method

Newton's method for finding the zero of a scalar function $f(x)$ is to iterate estimates of the zero and the local gradients, i.e.

$$x_{i+1} = x_i - \left[\frac{df(x_i)}{dx}\right]^{-1} f(x_i). \tag{8.109}$$

Figure 8.6 illustrates this: an initial estimate of the zero of $x_1$ together with the local gradient, $\frac{df(x_1)}{dx}$ is used to estimate $x_2$. The value of $x_2$ and the gradient at $x_2$ are then used to estimate $x_3$, and so on. For a vector valued function, $\mathbf{g(x)}$, the Newtonian iteration equation is

$$\mathbf{x}_{i+1} = \boldsymbol{\xi} - [\nabla\mathbf{g}(\boldsymbol{\xi})]^{-1} \mathbf{g}(\boldsymbol{\xi}). \tag{8.110}$$

### 8.6.2  Levenberg-Marquardt Method

The Levenberg-Marquardt method combines Newton's method with the method of steepest descent, using Newtonian iteration near the solution and steepest descent far from the solution. The Levenberg-Marquardt method can be used to find the minimum of $\chi^2$ from the algorithm

$$\mathbf{x}_{i+1} = \boldsymbol{\xi} - [\nabla_x^2\chi_i^2 - \gamma\mathbf{I}_n]^{-1}\nabla_x\chi_i^2 \tag{8.111}$$

where $\mathbf{I}_n$ is the $n \times n$ identity matrix. As $\gamma \to 0$, the Levenberg-Marquardt step tends to the Newtonian iteration method and as $\gamma \to \infty$ it tends to the method of steepest descent. The values of $\gamma$ are chosen to give steepest descent far from the solution and Newtonian iteration near the solution by monitoring $\chi^2$:

- If $\chi^2$ increases, $\gamma$ should be increased (making the step size smaller) and $\boldsymbol{\xi}$ left unchanged.

- If $\chi^2$ decreases, $\gamma$ should be decreased (making the step size larger) and $\boldsymbol{\xi}$ updated.

The factor by which $\gamma$ is increased or decreased can be flexible so it is somewhat dependent on the problem, but is most often set to a factor of 10.

Applying the Levenberg-Marquardt method to the optimal estimation problem where the cost function is given in equation 6.15 gives the iteration:

$$\mathbf{x}_{i+1} = \boldsymbol{\xi} + (\mathbf{S}_a^{-1} + \mathbf{K}_i^T \mathbf{S}_\epsilon^{-1} \mathbf{K}_i + \gamma \mathbf{I}_n)^{-1} [\mathbf{K}_i^T \mathbf{S}_\epsilon^{-1} (\mathbf{y} - \mathbf{F}(\boldsymbol{\xi})) - \mathbf{S}_a^{-1} (\boldsymbol{\xi} - \mathbf{x}_a)]. \quad (8.112)$$

As elements of the state vector may have different magnitudes and dimensions, they should be scaled. A convenient means of doing this is to replace $\mathbf{I}_n$ with $\mathbf{S}_a^{-1}$ as a scaling matrix, as it has the right dimensions. The Levenberg-Marquardt step then becomes:

$$\mathbf{x}_{i+1} = \boldsymbol{\xi} + (\mathbf{S}_a^{-1} + \mathbf{K}_i^T \mathbf{S}_\epsilon^{-1} \mathbf{K}_i + \gamma \mathbf{I}_n)^{-1} [\mathbf{K}_i^T \mathbf{S}_\epsilon^{-1} (\mathbf{y} - F(x_i)) - \mathbf{S}_a^{-1} (\boldsymbol{\xi} - \mathbf{x}_a)]. \quad (8.113)$$

Commonly in retrieval codes a sequence of iterations is deemed to have converged if the change in $\chi^2$ between steps is small. Alternatively, there is little point iterating if the change in $\mathbf{x}$ is much smaller than the uncertainty in $\mathbf{x}$ for all the components of $\mathbf{x}$.

### 8.6.3 The Non-Linear Solution

The parameters describing retrieval quality discussed in Section 8.5 (the averaging kernel, the degrees of freedom and the information content) were derived under the assumption that the forward model and retrieval method are linear. It is important to realise that they can still be applied for a non-linear retrieval provided the forward model is linear within the error bounds of the retrieval [*Rodgers*, 2000]. As the prior state appears only linearly in the iterative retrieval method, the same solution would be obtained by linearising the forward model at the retrieval solution. Thus, the averaging kernel is a valid concept for interpreting the iterative optimal estimation retrieval. Also, the total error covariance matrix of the solution $\hat{\mathbf{S}}$ is valid if it is evaluated at the solution and regarded as if the problem is linear at the solution.

A final important consideration in non-linear problems is that they may be ill-posed. Roughly speaking a problem is said to be well-posed if, for a wide class of data and in appropriate norms, it has a unique solution depending continuously on the data. In contrast, inverse problems are usually ill-posed and the $\chi^2$ surface may have multiple minima. Other minimization techniques may need to be employed (e.g. simulated annealing [*Press*, 1992]) to identify the region of the global minimum and so choose a suitable first guess from which the optimal estimation method can be employed.

## 8.7  Retrieval Error

The mapping from state space to measurement space has so far been assumed to be described by the forward model through

$$\mathbf{y} = \mathbf{F}(\mathbf{x}, \mathbf{b}) + \epsilon. \tag{8.114}$$

In many cases the measurement noise is the dominant error term, so $\mathbf{S}_\epsilon = \mathbf{S}_y$. However, in some cases significant forward model error compared with the physics is unavoidable. Two sources of error are considered

- the exclusion of some parameters in the forward model that may influence the behaviour of the forward function, e.g. in a model of transmission assuming an absorbing gas species can be ignored because it is usually present in very small amounts,

- the uncertainty in those parameters that are used in the forward model, e.g. in a model of transmission assuming a fixed amount of absorbing gas species when its amount actually varies.

The state and measurements are exactly related through the physical process embodied in the forward function

$$\mathbf{y} = \mathbf{f}(\mathbf{x}, \mathbf{b}) + \epsilon. \tag{8.115}$$

The forward model error, $\Delta\mathbf{f}$, is then characterised by

$$\Delta\mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') = \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') - \mathbf{F}(\mathbf{x}, \hat{\mathbf{b}}) \tag{8.116}$$

where $\mathbf{b}$ represents parameters used by the forward model and $\mathbf{b}'$ represents the forward function parameters which are ignored in the construction of the forward model. The term $\hat{\mathbf{b}}$ is used to denote the best estimates of $\mathbf{b}$ that are used in the forward model.

If the total measurement error (forward model errors plus experimental errors) are described by the total error covariance matrix $\mathbf{S}_\epsilon$ then the inverse solutions should be calculated using $\mathbf{S}_\epsilon$ rather than $\mathbf{S}_y$. However, it is worth examining the forward model error in detail.

### 8.7.1  Forward Model Error

The error analysis of the inverse method can be better understood if the problem is considered linear at the solution. If the forward model is then linearised about $\mathbf{x} = \mathbf{x}_a$ and $\mathbf{b} = \hat{\mathbf{b}}$ then

$$\mathbf{F}(\mathbf{x}, \mathbf{b}) = \mathbf{F}(\mathbf{x}_a, \hat{\mathbf{b}}) + \mathbf{K}_x(\mathbf{x} - \mathbf{x}_a) + \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) \tag{8.117}$$

where $\mathbf{K}_x = \frac{\partial \mathbf{F}}{\partial \mathbf{x}}$ and $\mathbf{K}_b = \frac{\partial \mathbf{F}}{\partial \mathbf{b}}$. The forward function can then be expressed as

$$\mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') = \mathbf{F}(\mathbf{x}, \hat{\mathbf{b}}) + \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') = \mathbf{F}(\mathbf{x}_a, \hat{\mathbf{b}}) + \mathbf{K}_x(\mathbf{x} - \mathbf{x}_a) + \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}').$$
(8.118)

The retrieved state $\hat{\mathbf{x}}$ comes from some inverse transform $\mathbf{R}$ that is a function of the measurements $\mathbf{y}$, forward model parameters $\hat{\mathbf{b}}$, the a priori knowledge of the state $\mathbf{x}_a$, and parameters that do not appear in the forward function but do affect the retrieval $\mathbf{c}$, i.e.

$$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{y}, \hat{\mathbf{b}}, \mathbf{x}_a, \mathbf{c}) \tag{8.119}$$

and the relation between the true state and the retrieved state becomes

$$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{f}(\mathbf{x}, \mathbf{b}) + \epsilon, \hat{\mathbf{b}}, \mathbf{x}_a, \mathbf{c}) \tag{8.120}$$

so that

$$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{F}(\mathbf{x}_a, \hat{\mathbf{b}}) + \mathbf{K}_x(\mathbf{x} - \mathbf{x}_a) + \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \epsilon, \hat{\mathbf{b}}, \mathbf{x}_a, \mathbf{c}). \tag{8.121}$$

Assuming the inverse method is linear with respect to $\mathbf{y}$ at the solution then

$$\hat{\mathbf{x}} = \mathbf{R}(\mathbf{F}(\mathbf{x}_a, \hat{\mathbf{b}}), \hat{\mathbf{b}}, \mathbf{x}_a, \mathbf{c}) + \mathbf{G}_y \left[ \mathbf{K}_x(\mathbf{x} - \mathbf{x}_a) + \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \epsilon \right] \tag{8.122}$$

where $\mathbf{G}_y = \frac{\partial \mathbf{R}}{\partial \mathbf{y}}$. It follows from the definition of $\mathbf{A}$ that

$$\mathbf{A} = \left( \frac{\partial \hat{\mathbf{x}}}{\partial \mathbf{x}} \right) = \mathbf{G}_y \mathbf{K}_x. \tag{8.123}$$

Following *Rodgers* [2000] an expression for the error in the retrieval $\hat{\mathbf{x}}$ as a difference from the true state can be obtained from Equation 8.122 by using the fact that

$$\mathbf{R}(\mathbf{F}(\mathbf{x}_a, \hat{\mathbf{b}}), \hat{\mathbf{b}}, \mathbf{x}_a, \mathbf{c}) = \mathbf{x}_a \tag{8.124}$$

so

$$\begin{aligned}
\hat{\mathbf{x}} - \mathbf{x} &= \mathbf{x}_a + \mathbf{G}_y \mathbf{K}_x(\mathbf{x} - \mathbf{x}_a) + \mathbf{G}_y \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \mathbf{G}_y \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \mathbf{G}_y \epsilon - \mathbf{x} \\
&= \mathbf{x}_a - \mathbf{x} + \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}_a + \mathbf{G}_y \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \mathbf{G}_y \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \mathbf{G}_y \epsilon \\
&= (\mathbf{A} - \mathbf{I}_n)\mathbf{x} - (\mathbf{A} - \mathbf{I}_n)\mathbf{x}_a + \mathbf{G}_y \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \mathbf{G}_y \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \mathbf{G}_y \epsilon \\
&= (\mathbf{A} - \mathbf{I}_n)(\mathbf{x} - \mathbf{x}_a) + \mathbf{G}_y \mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}}) + \mathbf{G}_y \Delta \mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}') + \mathbf{G}_y \epsilon. \tag{8.125}
\end{aligned}$$

The terms in this equation correspond to a source of error propagated through the inverse method into the solution:

- The first term, $(\mathbf{A} - \mathbf{I}_n)(\mathbf{x} - \mathbf{x}_a)$, is known as the smoothing error. This is the error due to the lack of sensitivity of the observing system to the individual parameters of the state vector. This term will be zero if on average $\mathbf{x} = \mathbf{x}_a$, i.e. the set of potential $\mathbf{x}$ is unbiased with respect to the a priori.

- The second term, $\mathbf{G}_y\mathbf{K}_b(\mathbf{b} - \hat{\mathbf{b}})$, is known as the model parameter error. The error covariance of this contribution is $\mathbf{G}_y\mathbf{K}_b\mathbf{S}_b\mathbf{K}_b^{\mathrm{T}}\mathbf{G}_y^{\mathrm{T}}$ where $\mathbf{S}_b = \langle(\mathbf{b} - \hat{\mathbf{b}})(\mathbf{b} - \hat{\mathbf{b}})^{\mathrm{T}}\rangle$. Typically $\mathbf{S}_b$ is a diagonal matrix with the elements of the diagonal being the uncertainties in the model parameters.

- The third term, $\mathbf{G}_y\Delta\mathbf{f}(\mathbf{x}, \mathbf{b}, \mathbf{b}')$, is known as the forward model error. If the forward model is based on a mathematical approximation then the forward model error is calculated as the typical difference between the approximation and the more exact model. In other cases, knowledge of the true physics may be so poor as to make estimates of the forward model error little more than an educated guess.

- The final term, $\mathbf{G}_y\boldsymbol{\epsilon}$, is known as the retrieval noise. It can be interpreted as the measurement noise projected into state space and its covariance is represented by $\mathbf{G}_y\mathbf{S}_y\mathbf{G}_y^{\mathrm{T}}$.

The full estimation of the uncertainty in the retrieval requires the evaluation and summation of the last three of these contributions. In practice, only the final term is usually considered. If a forward model parameter is so poorly known that it contributes significantly to the error budget, consideration should be given to including it in the state vector and so potentially estimating a more exact value. A target often used in developing a forward model is to require an accuracy that is about 1/10 th of the measurement uncertainty, thus ensuring the measurement error dominates the error budget.

---

**Problem 8.1** Show that if $\mathbf{S}_y = \mathbf{K}\mathbf{S}_x\mathbf{K}^{\mathrm{T}}$ then $\mathbf{S}_x^{-1} = \mathbf{K}^{\mathrm{T}}\mathbf{S}_y^{-1}\mathbf{K}$.

**Problem 8.2** Derive Equation 8.87 from Equation 8.86.

**Problem 8.3** Starting from Equation 8.89 derive the identity $\mathbf{G} = \mathbf{S}_a\mathbf{K}^{\mathrm{T}}(\mathbf{K}\mathbf{S}_a\mathbf{K}^{\mathrm{T}} + \mathbf{S}_y)^{-1}$.

**Problem 8.4** Derive Equation 8.91 from Equation 8.83.

**Problem 8.5** Derive the identity $\mathbf{I} - \mathbf{A} = \hat{\mathbf{S}}\mathbf{S}_a^{-1}$.

**Problem 8.6** Consider the retrieval example described in Section 8.3.2. Given an a priori value of the state of

$$\mathbf{x}_a = \begin{bmatrix} 101 \\ 1.8 \end{bmatrix}$$

which has an associated covariance matrix

$$\mathbf{S}_a = \begin{bmatrix} 4 & 0 \\ 0 & 0.5 \end{bmatrix}$$

calculate the following

1. the optimal solution,

2. the error covariance matrix for the optimal estimate,

3. the averaging kernel,

4. the degrees of freedom for signal and noise,

5. the change in information content caused by making the measurement.

**Problem 8.7** Generate the equation to estimate $\mathbf{x}_{i+1}$ using Newton's method by applying Equation 8.110 to the cost expression for maximum likelihood.

**Problem 8.8** A three band radiometer operating at 3.7, 10.8 and 12 $\mu m$ is used to observe the sea surface at night along an isothermal path. Measurements in the three channels can be assumed to be independent. The transmission of the atmosphere from the radiometer to the sea surface in each of the bands is $0.94 \pm 0.2$, $0.84 \pm 0.2$ and $0.95 \pm 0.2$ respectively. Similarly, the emissivity of sea water in each of the bands is 0.99, 0.98 and 0.96. A climatology of sea surface temperature suggests a value of $15 \pm 3$ C for the location observed by the radiometer. Given that the radiances measured by the instrument were $Z \pm z$, $X \pm x$ and $Y \pm y$, determine optimal estimates of the sea surface temperature and the intervening atmosphere and their uncertainty. Hint: this is a non-linear problem so is best solved using an iterative computer programme.

## Additional Reading

Rodgers, C. D., *Inverse Methods for Atmospheric Sounding: Theory and Practice*, World Scientific Pub Co Inc, Singapore, 2000.
Twomey, S., *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurements*, Dover Publications, New York, 1996