

Chapter 2

Mathematical Concepts

- Definition of measurement
- Mathematical formulae used to model *known* influences
- Normal distribution modeling *unknown* or *unquantifiable* influences
- Standard deviation is a measure for uncertainty

MATHEMATICAL concepts behind a measurement are described very formally in [2] the so-called GUM. This chapter rephrases some of the statements in the GUM, comments them and draws several consequences and implications on the automated system MUSAC.

2.1 What is a Measurement

In this section we summarize and comment the basic definitions, rules and assumptions on which the GUM [2, p. 31f.] is founded.

2.1.1 Definitions

The GUM defines several terms very precisely. Let's cite some of them here:

(measurable) quantity attribute of phenomenon, body or substance that may be distinguished qualitatively and determined quantitatively.

NOTES

- 1 The term quantity may refer to a **quantity in a general sense** [see example a)] or to a **particular quantity** [see example b)]

EXAMPLES

- a) quantities in a general sense: length, time, mass, temperature, electrical resistance, amount-of-substance concentration;
 - b) particular quantities:
 - length of a rod
 - electrical resistance of a given specimen of wire
 - amount-of-substance concentration of ethanol in a given sample of wine.
- 2 Quantities that can be placed in order of magnitude together into **categories of quantities**, for example:
 - work, heat, energy
 - thickness, circumference, wavelength.
 - 3 **Symbols for quantities** are given in ISO 31.

value (of a quantity) magnitude of a particular quantity generally expressed as a unit of measurement multiplied by a number [...]

NOTES

- 1 The value of a quantity may be positive, negative or zero.
- 2 The value of a quantity may be expressed in more than one way.
- 3 The values of quantities of dimension one are generally expressed as pure numbers.
- 4 A quantity that cannot be expressed as a unit of measurement multiplied by a number may be expressed by reference to a conventional reference scale or to a measurement procedure or to both.

measurement set of operations having the object of determining a value of a quantity.

NOTE: The operations may be performed automatically.

measurement procedure set of operations, described specifically, used in the performance of particular measurements according to a given method

NOTE: A measurement procedure is usually recorded in a document that is sometimes called a “measurement procedure” (or a **measurement method**) and is usually in sufficient detail to enable an operator to carry out a measurement without additional information.

measurand particular quantity subject to measurement

EXAMPLE: vapour pressure of given sample of water at 20 °C.

NOTE: The specification of a measurand may require statements about quantities such as time, temperature and pressure.

result of a measurement value attributed to a measurand, obtained by a measurement

NOTES

1 When a result is given, it should be made clear whether it refers to:

- the indication
- the uncorrected result
- the corrected result

and whether several values are averaged.

2 A complete statement of the result of a measurement includes information about the uncertainty of measurement.

uncertainty (of a measurement) parameter, associated with the result of a measurement, that characterizes the dispersion of the values that could reasonably be attributed to the measurand

NOTES

1 The parameter may be, for example, a standard deviation (or a given multiple of it), or the half-width of an interval having a stated level of confidence.

2 Uncertainty of measurement comprises, in general, many components. Some of these components may be evaluated from the statistical distribution of the result of series of measurements and can be characterized by experimental standard deviations. The other components which can also be characterized by standard deviations, are

evaluated from assumed probability distributions based on experience or other information.

- 3 It is understood that the result of the measurement is the best estimate of the value of the measurand, and that all components of uncertainty, including those arising from systematic effects, such as components associated with corrections and reference standards, contribute to the dispersion.

These definitions are not undisputed: E.g. Feller [20] and others criticize the definition of a measurable quantity having two different meanings: a quantity in a *general sense* is a very abstract notion of physical phenomena such as length, time, energy which are difficult to define themselves. Richard Feynman explained the concept of energy in [21]: "... there is a certain quantity, which we call energy, that does not change in the manifold changes which nature undergoes. That is a most abstract idea, because it is a mathematical principle... It is not a description of a mechanism, or anything concrete; it is just a strange fact that we can calculate some number and when we finish watching nature going through her tricks and calculate the number again, it is the same. ... It is important to realize that in physics today, we have no knowledge of what energy *is*." On the other hand, the term quantity may also refer to a *particular quantity* such as the "electrical resistance of a given specimen of wire" according to the metrological definition. This is a very different level of abstraction. Unfortunately, the observable phenomenon and its abstract principle are not clearly distinguished. This may sometimes lead to confusion: The length of two particular sticks are two different quantities which in turn are different to the quantity "length".

However, one is able to understand the crucial and central message of the GUM [2, § 3.1.2] with the definitions cited above:

In general, the result of a measurement is only an approximation or estimate of the value of the measurand and thus is complete only when accompanied by a statement of the uncertainty of that estimate.

In other words, the result of a measurement relates to an estimated mean and the uncertainty of the measurement refers to the estimated variance of a random variable. A comment on terminology: In this text the shorter term "measurement" is often used instead of the correct but clumsy "result of a measurement". The context will (hopefully) always clarify the meaning.

2.2 Different Approaches

Several methods are proposed in order to tackle the problem of *how to compute measurement uncertainty effectively and accurately*. They can be classified into two groups of different approaches. The one proposed by the GUM can roughly be summarized as *linearizing the problem* whereas the second approach can be subsumed as *simulating the problem*. The simulating approach is not explicitly mentioned in the first edition of the GUM and currently there are no plans for a second edition. But there are plans for supplemental guides. One of which will be concerned with the propagation of distributions. It will recommend the *Monte Carlo Simulation*, a computer-intensive process that samples many times from the input distributions, and, by evaluating the model for each sample, provides a distribution for the output [15]. This approach of simulating the input has become feasible only in modern times thanks to the tremendous improvement of computer speed.

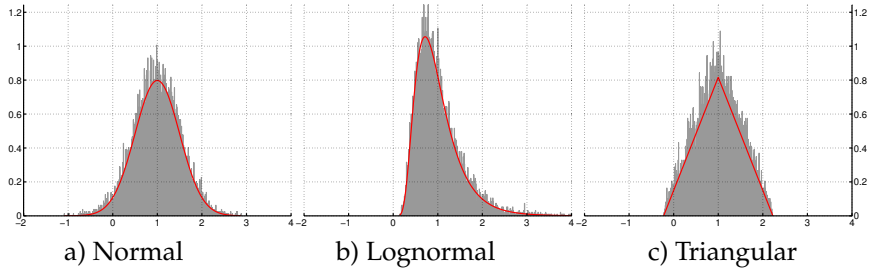
The following sections explain the reasoning and mathematics behind the linearizing approach, except section 2.2.5 which explains the Monte Carlo Simulation approach.

2.2.1 Random variables

Every single result of a measurement is unique and is subject to unique circumstances. Therefore no two executions of a measurement procedure yield the exact same result. One can imagine this fact as if infinitely many influences alter the result. This is why a result of a measurement can always be viewed as a *random variable*. A random variable is a statistical concept for describing situations such as the one above. Because there is no single value representing the “value of a random variable” one is interested in statistical statements on a random variable such as what is the *mean*, the *variance* or the *distribution* of its observed values. We will denote random variables by capital symbols such as X and Y . Instances of observed values of a random variable will be denoted by x_i or y_i . Conceptually, a random variable contains infinitely many observed values but for practical reasons there is often a finite number n of instances available. Think of a random variable as a basket containing n observed values.

Figure 1 shows different histograms of observed values. Observe that each of the three samples has the same mean $\mu = 1.0$ and the same vari-

Figure 1 Sampled random variables with different distributions. The mean and the variance are in each case the same $\mu = 1.0$ and $\sigma^2 = 0.25$.



ance $\sigma^2 = 0.25$. The diagrams are generated using 10000 observed random values which are depicted in a histogram of 500 columns. The height of the columns is normalized such that the area over all columns is 1. That is, as the number of observed values and the number of histogram columns grow big, these normalized histograms approach the *probability density function* of the underlying distribution which are also drawn in the histograms of Figure 1.

Mean The *mean* of a random variable X is a number describing the average of all observed values x_i . It is also called the *expected value* of X and it is denoted by $E[X] = \mu_X$. It is a single number that represents the average of infinitely many instances x_i . If we are observing n instances of X then the best we can do is to estimate the mean μ_X by

$$E[X] = \mu_X \approx \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (2.1)$$

With this definition the following properties of the expected value are obvious. They will be used in the following sections without any further reference:

$$E[\alpha] = \alpha \quad \text{and} \quad E[\alpha X + Y] = \alpha E[X] + E[Y],$$

where $\alpha \in \mathbb{R}$ and X and Y are arbitrary random variables.

Variance The *variance* $\text{Var}(X) = \sigma_X^2$ of a random variable X is a number describing the scattering of the observed values around the mean μ_X . The

larger the variance σ_X^2 is, the wider the instances are distributed around μ_X . The variance is the expected value of a random variable S defined as $S = (X - \mu_X)^2$. In other words: it is the average of the squared distance from the mean.

$$\sigma_X^2 = \text{Var}(X) = \text{E}[(X - \mu_X)^2] \quad (2.2)$$

Note the following important property:

$$\begin{aligned} \text{Var}(X) &= \text{E}[(X - \mu_X)^2] \\ &= \text{E}[X^2 - 2X\mu_X + \mu_X^2] \\ &= \text{E}[X^2] - 2\mu_X \text{E}[X] + \mu_X^2 \\ &= \text{E}[X^2] - 2\mu_X^2 + \mu_X^2 \\ \text{Var}(X) &= \text{E}[X^2] - \mu_X^2. \end{aligned} \quad (2.3)$$

If we are given n observations of X an estimate for σ_X^2 can be found by applying Equation (2.2) directly to the true mean μ_X if it is known: $\text{Var}(X) \approx 1/n \sum_{i=1}^n (x_i - \mu_X)^2$. But the true mean is almost never known, rather we can compute an estimate for the mean $\mu_X \approx \bar{X}$ using Equation (2.1). In that case the best estimate for the variance is found by:

$$\text{Var}(X) = \sigma_X^2 \approx s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2. \quad (2.4)$$

The division by $n - 1$ compensates for the fact that \bar{X} is an estimate rather than the true mean [17, 33]. Or one can apply Property (2.3). Compensating for the difference between true and estimated mean then again leads to

$$\text{Var}(X) = \sigma_X^2 \approx s_x^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right). \quad (2.5)$$

Covariance The *covariance* between two random variables X and Y expresses the relationship between the two of them. It is defined by

$$\text{Cov}(X, Y) = \text{E}[(X - \mu_X)(Y - \mu_Y)].$$

Observe that $\text{Cov}(X, X) = \text{Var}(X)$ and $\text{Cov}(X, Y) = \text{Cov}(Y, X)$. Independent random variables have a covariance equal to zero. The *correlation* $\rho_{X,Y}$ expresses the same relationship, but it is normalized such that it takes values between -1 and 1 :

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}. \quad (2.6)$$

The bigger the correlation's absolute value is, the stronger is the relationship between the two random variables. If the correlation, and therefore the covariance, is negative, then any two corresponding observed values are likely to lie on opposite sides of the respective mean. If the correlation is positive, then two corresponding observed values are probably on the same side of their mean.

For an estimation of the covariance we observe the following property:

$$E[(X - \mu_X)(Y - \mu_Y)] = E[XY - \mu_X \mu_Y] = E[XY] - \mu_X \mu_Y \quad (2.7)$$

This means that given two samples for two random variables of size n their covariance can be estimated either by

$$\text{Cov}(X, Y) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - n \bar{X} \bar{Y} \right) \quad (2.8)$$

or by

$$\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y}). \quad (2.9)$$

As before, the division by $n-1$ instead of by n is the best correction for the fact that we are using estimated instead of the true means.

There is, of course, a lot more to be said on random variables, statistics and so on. But I would rather point the interested reader to some good literature like [17, 28, 29, 33, 40].

Accurate computation of mean, variance and covariance In order to calculate the estimated mean a first approach would naïvely iterate over the data and add each summand to the partial sum $T_{i,j} = \sum_{k=i}^j x_k$ finally computing the total sum $T = T_{1,n}$. This means, especially for large i , that at step i adding x_{i+1} to $T_{1,i}$ eventually is like adding a drop of water to the sea. The relative error bound for $T = T_{1,n}$ computed in this way is basically $n\epsilon$,

n being the number of summands to add, and ε being a small error due to floating point representation if nothing else, the so called *machine epsilon*. In other words, the error for the sum grows with the number of summands involved. A good idea would be to sort the data and start summing up with the smallest numbers. The idea of this approach is to keep the intermediate sum small and thereby minimizing the absolute error per step. But still the partial sum eventually might grow so big that the addition of one single number would not affect the result, again due to round off. This is a fundamental problem when building a sum of many summands. William Kahan proposed Algorithm 1 to compute a sum of many summands. Its

Algorithm 1 Kahan's Summation Formula

```

function  $T = \text{Sum}(x, n)$ 
  — computation of  $T = \sum_{i=1}^n x_i$  —
   $T \leftarrow 0$ 
   $c \leftarrow 0$ 
  forall  $1 \leq i \leq n$  do
     $y \leftarrow x_i - c$            — add the previous cut off error —
     $T_{\text{old}} \leftarrow T$ 
     $T \leftarrow T + y$ 
     $c \leftarrow (T - T_{\text{old}}) - y$  — fetch the (negative) cut off error —
  end forall

```

idea is to store the round off error in each step and add it to the subsequent summand. This reduces the relative error bound for T to about 2ε which is a dramatic improvement compared to the naïve algorithm. See [26] and [30] for a thorough discussion of Kahan's summation formula.

In order to calculate the variance by Equation (2.4) the data must be traversed twice: once for finding \bar{X} and a second time for computing $\sum (x_i - \bar{X})^2$. In order to avoid these two passes statistical textbooks propose to reorganize it into Equation (2.5), building $\sum x_i$ and $\sum x_i^2$ in one pass. But this is deprecated, because of the possibly catastrophic numerical cancellation happening in the subtraction operation. Chan, Golub and LeVeque [14] propose to calculate the variance and mean in a recursive manner as shown in Algorithm 2. It is based on the so called *pairwise summation* — it deserves its name because pairs of similar magnitude are added at each recursion

level — and the *updating formula* by Youngs and Cramer [44]

$$\begin{aligned} T_{1,j} &= T_{1,j-1} + x_j \\ S_{1,j} &= S_{1,j-1} + \frac{1}{j(j-1)}(jx_j - T_{1,j})^2 \end{aligned} \quad (2.10)$$

with $T_{1,1} = x_1$ and $S_{1,1} = 0$ where $S_{ij} = \sum_{k=i}^j (x_k - 1/(j-i+1)T_{ij})^2$ is the sum of squares in analogy to T_{ij} as above. Chan, Golub and LeVeqe show that this approach is stable moreover they conjecture that the error bound is $\kappa \log(n)\varepsilon$ whereas Equation (2.4) is stable with an error bound $\kappa^2 n\varepsilon$ and Equation (2.5) is even unstable. Barlow [7] proves that the error bound of Algorithm 2 is indeed $\kappa \log(n)\varepsilon$ as conjectured, where $\kappa = \|x\|_2/\sqrt{5}$ is the condition number as defined in [14] for the calculation of the sample variance and $S = \sum (x_i - \bar{X})^2$.

In the paragraph above we showed very briefly the problem of calculating large sums i.e. the mean and the variance of a sample x along with some suggested solutions and refinements. In the context of MUSAC the goal is to have a fast and accurate algorithm, where memory and data availability is no problem. Several experiments showed that choosing the best algorithm is not obvious. Table 1 lists several criteria of various tested algorithms. The algorithms compared each perform the following four tasks:

1. find $x_{\min} = \min_i x_i$
2. find $x_{\max} = \max_i x_i$
3. compute $T = \sum_i x_i$
4. compute $S = \sum_i (x_i - \bar{X})^2$

The performance is given in milli seconds on a Pentium III 750 MHz. The two columns show different performances. The first column shows the time used if the executable code is not optimized by the compiler, the second column shows the time used, if all speed relevant optimizations of the compiler are turned on. Care must be taken not to “over-optimize” the executable code which can render Kahan’s summation formula futile. The last column shows the relative error of $|(S_{\text{sort}} - S)/S_{\text{sort}}|$, where S_{sort} is the exact result obtained by using a two pass algorithm sorting the data in each pass and using Kahan’s summation. S is the result of the examined algorithm. Finally, the

Algorithm 2 Recursive Calculation of Variance and Mean

```

function  $\langle m, v \rangle = \text{var}(x)$ 
  —  $x$ : sample  $\{x_i\}$ ,  $m$ :  $E[x]$ ,  $v$ :  $\text{Var}(x)$  —
   $n \leftarrow \text{length}(x)$ 
   $\langle T, S \rangle \leftarrow \text{var}_{\text{rec}}(x, 1, n)$  — call the recursive working horse —
   $m \leftarrow T/n$ 
   $v \leftarrow S/(n-1)$ 

function  $\langle T, S \rangle = \text{var}_{\text{rec}}(x, i, j)$ 
  —  $T$ :  $\sum_{k=i}^j x_k$ ,  $S$ : variance of subsample  $\sum_{k=i}^j (x_k - T/(j-i+1))^2$  —
  if  $i = j$  then
     $T \leftarrow x_i$ 
     $S \leftarrow 0$ 
  else
     $k \leftarrow \lfloor (i + j)/2 \rfloor$ 
     $m \leftarrow k - i + 1$ 
     $n \leftarrow j - k$ 
     $\langle T_1, S_1 \rangle \leftarrow \text{var}_{\text{rec}}(x, i, k)$ 
     $\langle T_2, S_2 \rangle \leftarrow \text{var}_{\text{rec}}(x, k + 1, j)$ 
     $T \leftarrow T_1 + T_2$ 
    if  $m = n$  then
       $S \leftarrow S_1 + S_2 + (T_1 - T_2)^2/2m$  — save a few flops —
    else
       $S \leftarrow S_1 + S_2 + \frac{m(n/m T_1 - T_2)^2}{n(m+n)}$ 
    end if
  end if

```

data for these experiments were a sample of 10^6 normally distributed random numbers with mean 1 and $\sigma = 10^{-8}$. Admittedly, this example is fairly extreme but it makes the differences more obvious. In more moderate situations (bigger σ) the relative error of S generally becomes smaller. The gist of

Table 1 Algorithms for computing a large sum

Algorithm	relative error bound		performance		relative error of S
	for T	for S	not	opt	
one pass, Equation (2.5)	$n\varepsilon$	$n\kappa^2\varepsilon$	451	130	∞
recursive pairwise summation	$\log(n)\varepsilon$	$\log(n)\kappa\varepsilon$	3184	932	3.1×10^{-11}
iterative pairwise summation	$\log(n)\varepsilon$	$\log(n)\kappa\varepsilon$	2153	501	8.6×10^{-12}
simple updating using (2.10)	$n\varepsilon$	$n\kappa\varepsilon$	470	361	1.5×10^{-9}
updating with Kahan's summation	2ε	$2\kappa\varepsilon$	581	450	7.6×10^{-12}
simple two pass	$n\varepsilon$	$n\varepsilon + n^2\kappa^2\varepsilon^2$	461	220	7.0×10^{-12}
two pass with Kahan's summation	2ε	$2\varepsilon + \mathcal{O}(\kappa^2\varepsilon^2)$	561	310	0
two pass, sorting, Kahan's summation	1ε	1ε	89058	6049	–

the experiments shown in Table 1 is that the two pass algorithm is not only the most accurate but also quicker than any other comparably accurate algorithm. It is quite noteworthy that the updating algorithms are slower than the two pass algorithms. This is attributed to the fact that the two pass algorithm uses two but particularly simple and thus very fast loops over the data rendering the two pass algorithm even faster than the updating one pass algorithm. Note, that the result for S obtained by the school book one pass algorithm according to Equation (2.5) was negative which is senseless for a

variance. Also note that the recursive and the iterative implementation of the *pairwise summation* yield slightly different results. This is due to the fact that the iterative implementation does not build the partial sums in *exactly* the same manner as the recursive implementation does. But it is very obvious that the overhead needed for storing and managing the partial sums is quite expensive. Therefore, Algorithm 3 is used in the MUSAC-system for finding the four characteristic numbers as stated on page 18 in Monte Carlo Simulations.

Algorithm 3 Mean, Variance, Minimum and Maximum of a Sample

```

function  $\langle m, v, a, b \rangle = \text{update}(x)$ 
  —  $m$ : mean,  $v$ : variance,  $a$ : minimum,  $b$ : maximum of sample  $x$  —
   $a \leftarrow x_1, \quad b \leftarrow x_1, \quad T \leftarrow 0, \quad S \leftarrow 0$  — initializations —
   $c \leftarrow 0, \quad n \leftarrow \text{length}(x)$  — temporary variables —
  forall  $1 \leq i \leq n$  do
     $y \leftarrow x_i$ 
     $a \leftarrow \min(a, y), \quad b \leftarrow \max(b, y)$ 
     $y \leftarrow y - c$  — Kahan's summation —
     $T_{\text{old}} \leftarrow T$ 
     $T \leftarrow T + y$ 
     $c \leftarrow (T - T_{\text{old}}) - y$  — do not "optimize" this statement —
  end forall
   $m = T/n$ 
   $c \leftarrow 0$ 
  forall  $1 \leq i \leq n$  do
     $y \leftarrow x_i - m, \quad y \leftarrow y^2$  — simple and quick instructions —
     $y \leftarrow y - c$  — Kahan's summation again —
     $T_{\text{old}} \leftarrow T$ 
     $T \leftarrow T + y$ 
     $c \leftarrow (T - T_{\text{old}}) - y$ 
  end forall
   $v \leftarrow S/(n-1)$  — c.f. Equation (2.4) —

```

For computing the covariance Equation (2.8) is numerically unstable for the same reason that Equation (2.5) is unstable whereas Equation (2.9) is stable. The experiments for a quick and accurate computation of the variance

showed that the two pass algorithm is not only very accurate but also comfortably quick. Thus we apply the same technique to the calculation of the covariance. So in the first pass we calculate the means of the two samples.

Note, that since the means are calculated in a first pass, we could shift the data by their respective mean rendering the body of the double loop as simple as $S \leftarrow x_i y_i$. But experiments showed that the access to the additional memory for storing the shifted data is by far more time consuming than executing the subtraction operation within the loop. Thus Algorithm 4 turned out to be the fastest as well as accurate for computing the covariance.

Algorithm 4 Calculation of the Covariance

```

function  $c = \text{cov}(x, y)$ 
   $c \leftarrow 0, \quad n \leftarrow \text{length}(x)$       — initialization —
   $\bar{X} \leftarrow \text{mean}(x)$                     — first pass to find means —
   $\bar{Y} \leftarrow \text{mean}(y)$ 
  forall  $1 \leq i < n$  do
     $c \leftarrow c + (x_i - \bar{X})(y_i - \bar{Y})$  — using Kahan's summation once more —
  end forall
   $c \leftarrow c/n-1$ 

```

2.2.2 Mathematical Model

To repeat the execution of a measurement n times is in most cases infeasible and often it is even impossible. Additionally, the result of a measurement Y most often depends on many other quantities X_i . The GUM assumes that this dependence can be modeled using mathematical functions. Thus any measurement can be expressed through a functional relationship

$$Y = f(X_1, X_2, \dots, X_n), \quad (2.11)$$

where the X_i are all quantities that influence the measurement Y . It is important that the effects of the influences X_i are known or at least can be modeled by a mathematical function as well.

Law of Propagation of Uncertainty Consider a quantity Y depending on n influences as in Equation (2.11). We expand f around the expected value of X_i i.e. $E[X_i] = \mu_i$ in a Taylor series. This yields

$$Y = f(\mu_1, \mu_2, \dots, \mu_n) + \sum_{i=1}^n \frac{\partial f}{\partial X_i} (X_i - \mu_i) + \mathcal{O}\left(\sum (X_i - \mu_i)^2\right), \quad (2.12)$$

and for

$$\begin{aligned} E[Y] &\approx E\left[f(\mu_1, \mu_2, \dots, \mu_n) + \sum_{i=1}^n \frac{\partial f}{\partial X_i} (X_i - \mu_i)\right] \\ &= E[f(\mu_1, \mu_2, \dots, \mu_n)] + \sum_{i=1}^n \frac{\partial f}{\partial X_i} E[(X_i - \mu_i)] \\ &= f(\mu_1, \mu_2, \dots, \mu_n), \end{aligned}$$

where $\partial f/\partial X_i$ is the derivative of f with respect to X_i evaluated at $\mu_1, \mu_2, \dots, \mu_n$. The GUM assumes that higher order terms are mostly negligible. The first order approximation for $E[Y]$ is

$$E[Y] = \mu_Y \approx \bar{Y} = f(\mu_1, \mu_2, \dots, \mu_n). \quad (2.13)$$

In Equation (2.12) we cut the Taylor series expansion after the linear term. The next following term of the expansion is:

$$\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial X_i \partial X_j} (X_i - \mu_i)(X_j - \mu_j) + \mathcal{O}\left(\sum (X_i - \mu_i)^3\right). \quad (2.14)$$

We now get an additional term which in turn will lead to a better approximation for the expected value of $E[Y]$ as follows:

$$\begin{aligned} E[Y] &= \mu_Y \\ &\approx \bar{Y} = E[f(\mu_1, \mu_2, \dots, \mu_n)] + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial X_i \partial X_j} E[(X_i - \mu_i)(X_j - \mu_j)] \\ &= f(\mu_1, \mu_2, \dots, \mu_n) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 f}{\partial X_i \partial X_j} \text{Cov}(X_i, X_j) \end{aligned} \quad (2.15)$$

Thus we get a second order approximation for the expectation of Y . The additional term of Equation (2.15) may be split into two sums. One sum is

taken over the variances and the other is taken over the covariances of X_i , this leads to the following expression for the second order approximation:

$$E[Y] \approx f(\mu_1, \mu_2, \dots, \mu_n) + \frac{1}{2} \sum_{i=1}^n \frac{\partial^2 f}{\partial X_i^2} \text{Var}(X_i) + \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{\partial^2 f}{\partial X_i \partial X_j} \text{Cov}(X_i, X_j).$$

In the same way we need to estimate $\text{Var}(Y)$ the variance of Y , with Y again given as in Equation (2.12). Substituting the Taylor Series Expansion of f (2.12) in Property (2.3) yields

$$\begin{aligned} \text{Var}(Y) &= E\left[(Y - \mu_Y)^2\right] \\ &\approx E\left[\left(f(\mu_1, \mu_2, \dots, \mu_n) + \sum_{i=1}^n \frac{\partial f}{\partial X_i} (X_i - \mu_i) - \mu_Y\right)^2\right] \\ &\approx \sum_{i=1}^n \sum_{j=1}^n \frac{\partial f}{\partial X_i} \frac{\partial f}{\partial X_j} E[(X_i - \mu_i)(X_j - \mu_j)], \end{aligned}$$

which — after splitting the sum into variance and covariance terms — leads to the following approximation for $\text{Var}(Y)$ as stated in [2, §5.2.2]:

$$\text{Var}(Y) \approx \sum_{i=1}^n \left(\frac{\partial f}{\partial X_i}\right)^2 \text{Var}(X_i) + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{\partial f}{\partial X_i} \frac{\partial f}{\partial X_j} \text{Cov}(X_i, X_j). \quad (2.16)$$

Expanding the Taylor series of Y^2 up to quadratic order leads to (now using the quadratic order approximation (2.15) for $E[Y]$) the following additional term to be added to Equation (2.16).

$$\sum_{i=1}^n \sum_{j=1}^n \left(\frac{1}{2} \left(\frac{\partial^2 f}{\partial X_i \partial X_j}\right)^2 + \frac{\partial f}{\partial X_i} \frac{\partial^3 f}{\partial X_i \partial X_i \partial X_j^2}\right) \text{Var}(X_i) \text{Var}(X_j). \quad (2.17)$$

The GUM assumes that first order terms are in most cases sufficient. But it is not entirely clear whether or not higher order terms should be used. On page 19, § 5.1.2 the GUM [2] reads: “When the nonlinearity of f is significant, higher-order terms in the Taylor series expansion must be included in the expression for $\text{Var}(Y)$...”. Unfortunately the GUM does not say what exactly a *significant* nonlinearity of f really is. But it states, that “... when the distribution of X_i is symmetric about its mean, the most im-

portant terms of next highest [higher] order to be added to terms of Equation (2.16) are (2.17)". Under the constraints of symmetric distributions the second order estimate for $\text{Var}(Y)$ is therefor the sum of (2.16) and (2.17).

Occasionally, one is interested in the covariance of two random variables whose relation is known. Suppose we are looking for the covariance of two random variables that are constructed by different combinations of one random variable, i.e. $\text{Cov}(f(X), g(X))$, or we might be asking for the covariance between the random variables X and $Z = h(X, Y)$ which is constructed using X and Y with a known covariance $\text{Cov}(X, Y) = c$. In the next two paragraphs we will derive expressions for the two cases in order to show how these questions can be answered mechanically. For simplicity we will stick to the first order approximations but expressions for the second order approximation can be found analogously. As a first step, let us write down the first order approximations of the functions involved:

$$F = f(X) \approx \bar{f} + \frac{df}{dX}(X - \bar{X})$$

$$G = g(X) \approx \bar{g} + \frac{dg}{dX}(X - \bar{X})$$

$$H = h(X, Y) \approx \bar{h} + \frac{\partial h}{\partial X}(X - \bar{X}) + \frac{\partial h}{\partial Y}(Y - \bar{Y}),$$

where $\bar{f} = f(\bar{X})$ and so forth. In the first example we are looking for

$$\text{Cov}(F, G) = E[FG] - E[F]E[G]$$

$$\begin{aligned} &\approx E\left[\left(\bar{f} + \frac{df}{dX}(X - \bar{X})\right)\left(\bar{g} + \frac{dg}{dX}(X - \bar{X})\right)\right] - \bar{f}\bar{g} \\ &= E\left[\bar{f}\bar{g} + \bar{g}\frac{df}{dX}(X - \bar{X}) + \bar{f}\frac{dg}{dX}(X - \bar{X}) + \frac{df}{dX}\frac{dg}{dX}(X - \bar{X})^2\right] - \bar{f}\bar{g} \\ &= \bar{f}\bar{g} + \left(\bar{f}\frac{dg}{dX} + \bar{g}\frac{df}{dX}\right)E[(X - \bar{X})] + \frac{df}{dX}\frac{dg}{dX}E[(X - \bar{X})^2] - \bar{f}\bar{g} \\ &= \frac{df}{dX}\frac{dg}{dX}\text{Var}(X). \end{aligned}$$

We only used Property (2.3) and the fact that $E[(X - \bar{X})] = 0$. In the second

example we make additional use of Property (2.7) to derive

$$\begin{aligned}
 \text{Cov}(X, H) &= E[XH] - E[X]E[H] \\
 &\approx E\left[X\left(\bar{h} + \frac{\partial h}{\partial x}(X - \bar{X}) + \frac{\partial h}{\partial y}(Y - \bar{Y})\right)\right] - \bar{X}\bar{h} \\
 &= E\left[\bar{h}X + \frac{\partial h}{\partial x}X(X - \bar{X}) + \frac{\partial h}{\partial y}X(Y - \bar{Y})\right] - \bar{X}\bar{h} \\
 &= \bar{X}\bar{h} + \frac{\partial h}{\partial x}E[X(X - \bar{X})] + \frac{\partial h}{\partial y}E[X(Y - \bar{Y})] - \bar{X}\bar{h} \\
 &= \frac{\partial h}{\partial x}\left(E[X^2] - \bar{X}^2\right) + \frac{\partial h}{\partial y}\left(E[XY] - \bar{X}\bar{Y}\right) \\
 &= \frac{\partial h}{\partial x}\text{Var}(X) + \frac{\partial h}{\partial y}\text{Cov}(X, Y).
 \end{aligned}$$

In the general case the following holds

$$\text{Cov}(f(x_1, \dots, x_n), g(y_1, \dots, y_m)) \approx \sum_{i=1}^n \sum_{j=1}^m \frac{\partial f}{\partial x_i} \frac{\partial g}{\partial y_j} \text{Cov}(x_i, y_j). \quad (2.18)$$

Fundamental Steps Equipped with this technique the general procedure for evaluating uncertainty according to the GUM is obvious:

- Build a mathematical model for the measurement using all possible (i.e. quantifiable) influences on the result.
- Determine the uncertainty of the result by applying the law of uncertainty propagation as shown in the previous sections.

The EURACHEM Guide [18, p. 12] expands them into the following four steps to be executed for every measurement:

1. Specify the measurement
2. Identify uncertainty sources
3. Quantify all influences
4. Determine the result of measurement and the measurement uncertainty

2.2.3 ISO Guide: First Order Approximation

The GUM proposes to apply the *law of propagation of uncertainty* in order to determine the measurement uncertainty of any measurement. It is based on various assumptions on the distributions of the random variables. Furthermore it assumes that the functions involved are sufficiently linear, that is, it assumes that nonlinearities can safely be neglected. Once these assumptions are accepted, the process for estimating the measurement uncertainty is indeed straightforward. Given the function

$$Y = f(X_1, X_2, \dots, X_n),$$

one applies Equation (2.14) for the estimation of $E[Y]$ and Equation (2.16) for the estimation of the variance $\text{Var}(Y)$ both up to first order. This process is quite a tedious and mechanical task, even more so in the case of second order approximations, c.f. Equations (2.15) and (2.17). Still, it should be mentioned that apart from these simplifying assumptions the main problem is the process of constructing an appropriate function $f(\cdot)$ that describes the physical laws of the measurement to the level of detail needed for the purpose of the measurement.

2.2.4 Simple Rules

After considering a variety of measurements it turns out that there is often no need to calculate the derivatives of the measurement function explicitly. The GUM [2, § 5.1.3, § 5.1.6] proposes the application of the following simple rules:

$$\begin{array}{ll} \text{if } Z = X \pm Y & \text{then } \begin{cases} E[Z] \leftarrow E[X] \pm E[Y] \\ \text{Var}(Z) \leftarrow \text{Var}(X) + \text{Var}(Y) \end{cases} \\ \text{if } Z = X \cdot Y & \text{then } \begin{cases} E[Z] \leftarrow E[X] \cdot E[Y] \\ \text{Var}(Z) \leftarrow \bar{Y}^2 \text{Var}(X) + \bar{X}^2 \text{Var}(Y) \end{cases} \end{array}$$

In other words, the variance and the mean can be calculated almost simultaneously. *But* these rules are only correct in the sense of first order approximation if

1. there are no covariances between the input variables, *and*
2. the input variables affect the measurement in exactly one distinct way only.

As an example for the second point, consider the case of $Z = X \cdot X$. The simple rules would not take into account the correlation between X and X which happens to be as strong a correlation as can be. As a rule of thumb these very simple rules cannot be applied wherever one input variable appears at more than once in the model formula and it cannot be applied if covariances between input variables are considered. And finally these simple rules cannot be applied when functions such as $\exp(\cdot)$, $\log(\cdot)$, $\sqrt{\cdot}$, ... are used.

2.2.5 Monte Carlo Simulation

In some cases the assumptions of Section 2.2.3 are too restrictive. In such cases other techniques must be used to reduce the errors induced by the assumptions. One way would be to employ higher order approximations. This approach would reduce errors due to nonlinearity in the model function $f(\cdot)$, but it still relies on the fact that distributions of the input variables would be symmetric. It turns out that most measurements are deliberately tuned to be as linear as possible, so that the need for higher order approximation due to nonlinearities is practically nonexistent. The reason for unacceptable errors produced by the method of Section 2.2.3 are due to violations of the assumed symmetric distributions. They appear in situations where a variable refers to a number for example close to 0% or close to 100%. E.g. it does not make sense to talk about a "normally distributed concentration with an average of 0.1% and a variance of 0.05%" because this would imply that with some considerable and non vanishing probability the "observed values" of that concentration are negative which is physical nonsense. Still, it might be the case that the concentration has indeed an observed mean of 0.1% and an observed variance of 0.05%, but the distribution of that concentration cannot be normal. By physical common sense it must be somehow skewed; as a matter of fact the most probable value for the concentration practically never coincides with the expected value as is

visualized in Figure 1 on page 14.

The expression *Monte Carlo Simulation* stands for *simulating true values*. The Monte Carlo Simulation is as simple as effective. It is based on the assumption that each *observed value* of a random variable can be treated as if it were the *true value*. The simulation therefore proceeds in the following steps:

- Generate n random values for each input variable
- Execute the model function $f(\dots)$ for each observed value assuming it to be the “true value”, thus generating n “true values” for Y which according to this assumption are interpreted as observed values of Y .

Now, in order to estimate the mean $E[Y]$ and the variance $\text{Var}(Y)$ Equations (2.1) and (2.4) can be applied. This method is indeed very simple in principle.

But let’s see how big n must be in order to yield reasonably accurate estimates! The expression for \bar{X} in Equation (2.1) is itself a random variable since we pick at random n values out of the infinitely many observable values of X . The variance of the random variable \bar{X} is proportional to $1/\sqrt{n}$ as explained in any handbook on statistics [19]. This means increasing the number of observed values n by a factor of 100 decreases the expected difference $\bar{X} - \mu_X$ by a factor of 10. This factor is equally bad for the estimation of the variance since the variance of the variance estimated according to Equation (2.4) is proportional to $1/\sqrt{n}$ as well, see [17]. Thus, the number n of observed values needed quickly grows very big. Even for a modest requirement of three correct digits one needs to generate observed values in the order of 10’000’000 for each input variable. Executing this method by hand is beyond discussion, but executing the same function n times by a computer is very attractive, since it is very simple to implement and computers are becoming ever faster. Therefore waiting is ever less an issue. Despite the simplicity of the principle of Monte Carlo simulations, it is very interesting to implement the handling of such large random variable samples efficiently. Section 5.3.3 on page 118 describes the techniques applied in some more details.

2.3 Calibration

This section was developed together with Oscar Chinellato [6]. I thank him for his time, his enthusiasm and the fun we had.

The following sections apply the principle of Section 2.2.2 thoroughly to the *calibration process* or a *general regression*. These are metrological and mathematical terms for the same thing: fitting a curve through some given points.

2.3.1 Calibration and Measurement

Most analytical methods are relative to some references and therefore need calibration. Thus calibration measurements are normally performed based on *reference material* or *calibration standards*. Usually least-squares methods ignoring the uncertainties associated with calibration standards are used. But different authors [12, 5] have shown that taking into account the uncertainties associated with calibration standards leads to better approximations of the model. In [12, 5] this is called a *Maximum Likelihood (fitting of a) Functional Relationship (MLFR)*. It is also shown in a terse way how the computation of MLFR can be executed. We present here two slightly different approaches. The first method, the so-called XIP-FIT, is very similar to the MLFR but it allows to take covariances between input data into account. The second method, the so-called P-FIT, yields results — at least in our tests — close to the *effective variance* approach as described in [39], but again it allows to take covariances among input data into account.

Since measurement instruments obey physical laws, they can be modeled by mathematical functions. Unless all free parameters \mathbf{p} are quantified these functions describe only the qualitative behavior. Calibration is the process of quantifying these parameters. For example, consider the measurement of a weight using a spring balance.

Here we denote weights by $x \in \mathbb{R}$ and stretches of the spring by $y \in \mathbb{R}$. Let x_{cs} be a known weight of a calibration standard and y_{cs} be the stretch of the spring balance for this calibration standard. For a given calibrated balance \mathbf{p} the following holds: $y_{cs} = f(x_{cs}, \mathbf{p})$. The function f is called *calibration function*. The art of calibrating consists of adjusting the parameters \mathbf{p} (e.g. the spring constant) using several calibrating pairs x_i and y_i such that $y_i \approx f(x_i, \mathbf{p})$ holds for all i as well as possible.

For an unknown mass \hat{x} , *weighing with a calibrated spring balance* requires the observation of a spring stretch \hat{y} and subsequent application of the *measurement function* $g(\hat{y}, \mathbf{p})$. In other words, $\hat{x} \approx g(\hat{y}, \mathbf{p})$. Note that g is the inverse function of f .

In the case that the pairs (x_i, y_i) were exact this would be a standard problem. But in practice every measurement is subject to uncertainties. Therefore, not only the calibrating pairs (x_i, y_i) but also their respective uncertainties (α_i, β_i) must be known and taken into consideration. Throughout this section the error random variables $e_i^{(x)}$ and $e_i^{(y)}$ are treated as $e_i^{(x)} \sim \mathcal{N}(0, \alpha_i)$ and $e_i^{(y)} \sim \mathcal{N}(0, \beta_i)$ respectively. They represent the distribution of the true values. This is in full compliance with the “Guide to the Expression of Uncertainty in Measurement” (GUM) [2].

With respect to our spring balance example, this implies that the real calibration weights are not x_i but $x_i + e_i^{(x)}$. Analogously the correct values for the stretch are no longer y_i but $y_i + e_i^{(y)}$.

Note that the uncertainties of x and y are uncorrelated. That is, the error of the calibration standard is independent of the error of the reading.

2.3.2 Calibration Models

Before we start considering calibration models, it is necessary to agree on a number of assumptions.

The measurement tool is expected to work according to the function $f(\cdot, \cdot)$ in use. This means that there exists a set of parameters which describes the apparatus correctly.

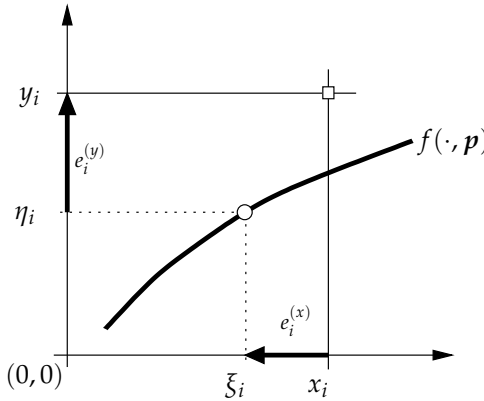
Every error *random vector* \mathbf{u} is assumed to be normally distributed and is denoted by $\mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ where $\boldsymbol{\mu} \in \mathbb{R}^n$ is the *mean vector* and $\mathbf{C} \in \mathbb{R}^{n \times n}$ is the *covariance matrix* of the random variable $\mathbf{u} \in \mathbb{R}^n$. The *multivariate probability density function* (pdf) is

$$\text{pdf}(\mathbf{u}) = \text{err}(\mathbf{u} - \boldsymbol{\mu}, \mathbf{C}) = (2\pi)^{-\frac{n}{2}} \frac{1}{\sqrt{\det(\mathbf{C})}} \exp\left(-\frac{1}{2}(\mathbf{u} - \boldsymbol{\mu})^\top \mathbf{C}^{-1}(\mathbf{u} - \boldsymbol{\mu})\right).$$

A crucial point for the further development of calibration models is the exact specification of the calibration process. A reasonable assumption for one calibration measurement is the following (the process is schematically depicted in Figure 2):

- Choose the i -th calibration standard x_i (whose real but unknown value ξ_i differs from x_i by $e_i^{(x)}$).
- Let \mathbf{p} be the real but unknown set of parameters that describes the tool exactly. Then η_i shall denote the true but (again) unknown value $f(\xi_i, \mathbf{p})$.
- Finally, the reading mechanism of the tool generally imposes an additional error $e_i^{(y)}$. This leads to $y_i = \eta_i + e_i^{(y)}$.

Figure 2 Process of a Calibration Measurement



Berkson's Model The *Berkson Model* [9, 13, 22, 31] is a statistical model that handles exactly the process described above and shown in Figure 2. It behaves as follows

$$\begin{aligned}\xi_i &= x_i + e_i^{(x)} \\ y_i &= f(\xi_i, \mathbf{p}) + e_i^{(y)}\end{aligned}\tag{2.19}$$

where x_i is a known and constant value and the errors $e_i^{(x)}$ and $e_i^{(y)}$ are independent and normally distributed with means 0 and variances α_i^2 and β_i^2 respectively.

A *Taylor series expansion* of $f(\xi_i, \mathbf{p})$ around x_i yields an expression that reveals the connection between ξ_i and y_i

$$y_i \approx f(x_i, \mathbf{p}) + \frac{\partial}{\partial x} f(x_i, \mathbf{p}) \cdot e_i^{(x)} + e_i^{(y)} = f^{(i)} + f_x^{(i)} e_i^{(x)} + e_i^{(y)}, \quad (2.20)$$

where $f^{(i)} = f(x_i, \mathbf{p})$ and $f_x^{(i)} = \frac{\partial}{\partial x} f(x_i, \mathbf{p})$. Equations (2.19) and (2.20) make the following relations obvious:

$$\begin{aligned} \text{Var}(\xi_i) &= \text{Var}(x_i + e_i^{(x)}) &&= \alpha_i^2 \\ \text{Var}(y_i) &\approx \text{Var}\left(f^{(i)} + f_x^{(i)} e_i^{(x)} + e_i^{(y)}\right) &&= f_x^{(i)2} \alpha_i^2 + \beta_i^2 \\ \text{Cov}(\xi_i, y_i) &\approx \text{E}\left[e_i^{(x)} \left(f_x^{(i)} e_i^{(x)} + e_i^{(y)}\right)\right] &&= f_x^{(i)} \alpha_i^2. \end{aligned}$$

Fitting of ξ and \mathbf{p} (XIP-FIT) As can be seen in Section 2.3.2, the pdf of (ξ_i, y_i) can be written as a *bivariate normal distribution*

$$\text{pdf}(\xi_i, y_i) = \text{err}(z_i, C_i) = \frac{1}{2\pi} \frac{1}{\sqrt{\det(C_i)}} \exp\left(-\frac{1}{2} z_i^\top C_i^{-1} z_i\right),$$

where z_i and C_i are defined as follows:

$$z_i = \begin{pmatrix} \xi_i - x_i \\ y_i - f^{(i)} \end{pmatrix} \quad \text{and} \quad C_i = \begin{pmatrix} \alpha_i^2 & f_x^{(i)} \alpha_i^2 \\ f_x^{(i)} \alpha_i^2 & f_x^{(i)2} \alpha_i^2 + \beta_i^2 \end{pmatrix}.$$

We generalize to all n measurements, and allow additional correlations among the errors $e_i^{(x)}$ and $e_i^{(y)}$ denoted by C_x and C_y respectively. This leads to the

$$\text{pdf}(\xi, \mathbf{y}) = \text{err}(z, C) = (2\pi)^{-n} \frac{1}{\sqrt{\det(C)}} \exp\left(-\frac{1}{2} z^\top C^{-1} z\right),$$

where z and C now are defined as

$$z = \begin{pmatrix} \xi - x \\ \mathbf{y} - \mathbf{f} \end{pmatrix} \quad \text{and} \quad C = \begin{pmatrix} C_x & C_x F_x \\ F_x C_x & F_x C_x F_x + C_y \end{pmatrix},$$

and analogously,

$$\mathbf{f} = \left(f^{(1)}, \dots, f^{(n)}\right)^\top \quad \text{and} \quad F_x = \text{diag}\left(f_x^{(1)}, \dots, f_x^{(n)}\right).$$

Thus we know the value of the pdf at (ξ, \boldsymbol{p}) given \boldsymbol{x} and \boldsymbol{y} . The maximum likelihood principle aims at maximizing this value. Hence

$$\max_{\xi, \boldsymbol{p}} \left((2\pi)^{-n} \frac{1}{\sqrt{\det(\boldsymbol{C})}} \exp\left(-\frac{1}{2} \boldsymbol{z}^\top \boldsymbol{C}^{-1} \boldsymbol{z}\right) \right) \quad (2.21)$$

must be solved for \boldsymbol{p} and ξ . See Section 2.3.3 for computational details.

Fitting of \boldsymbol{p} (P-FIT) In the previous section we fitted for both parameters \boldsymbol{p} and ξ . Hence we found the most probable pairing of \boldsymbol{p} and ξ , i.e. an estimation for the true shape of f and simultaneously the true values of the calibration standards \boldsymbol{x} . However, it is often the case that one is not interested in an estimate for the calibration standard values, but rather one is interested in the most probable shape of the calibration function f regardless of any estimate for \boldsymbol{x} .

Consequently, we are looking for the pdf of \boldsymbol{y} only. This is computed by the *law of total probability* [38]. We show the derivation of it for one calibration pair (x, y) . We use $f = f(x, \boldsymbol{p})$ and $f_x = \frac{\partial}{\partial x} f(x, \boldsymbol{p})$ as abbreviations.

$$\begin{aligned} \text{pdf}(\boldsymbol{y}) &= \int_{-\infty}^{\infty} \text{pdf}(\boldsymbol{y} | \xi) \text{pdf}(\xi) d\xi \\ &= \int_{-\infty}^{\infty} \text{err}(\boldsymbol{y} - f(\xi, \boldsymbol{p}), \beta) \text{err}(\xi - x, \alpha) d\xi \\ &\approx \int_{-\infty}^{\infty} \text{err}(\boldsymbol{y} - f - f_x(\xi - x), \beta) \text{err}(\xi - x, \alpha) d\xi \\ &= \frac{1}{2\pi} \frac{1}{\alpha\beta} \int_{-\infty}^{\infty} \exp\left(-\frac{1}{2} \frac{(\boldsymbol{y} - f - f_x(\xi - x))^2}{\beta^2}\right) \exp\left(-\frac{1}{2} \frac{(\xi - x)^2}{\alpha^2}\right) d\xi \\ &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{\alpha^2 f_x^2 + \beta^2}} \exp\left(-\frac{1}{2} \frac{(\boldsymbol{y} - f)^2}{\alpha^2 f_x^2 + \beta^2}\right) \\ \implies \text{pdf}(\boldsymbol{y}) &\approx \text{err}\left(\boldsymbol{y} - f, \sqrt{\alpha^2 f_x^2 + \beta^2}\right) \end{aligned}$$

The generalization to the n -dimensional case is now straightforward:

$$\text{pdf}(\mathbf{y}) \approx \text{err}(\mathbf{z}, \mathbf{C}) = (2\pi)^{-\frac{n}{2}} \frac{1}{\sqrt{\det(\mathbf{C})}} \exp\left(-\frac{1}{2}\mathbf{z}^\top \mathbf{C}^{-1}\mathbf{z}\right)$$

with

$$\mathbf{z} = \mathbf{y} - \mathbf{f} \quad \text{and} \quad \mathbf{C} = \mathbf{F}_x \mathbf{C}_x \mathbf{F}_x + \mathbf{C}_y.$$

Hence we know the value of the pdf at \mathbf{p} given \mathbf{x} and \mathbf{y} . Again using the principle of maximum likelihood, we are left with the maximization problem

$$\max_{\mathbf{p}} \left((2\pi)^{-\frac{n}{2}} \frac{1}{\sqrt{\det(\mathbf{C})}} \exp\left(-\frac{1}{2}\mathbf{z}^\top \mathbf{C}^{-1}\mathbf{z}\right) \right). \quad (2.22)$$

Superficially, this looks equivalent to Equation (2.21), but recall that \mathbf{C} and \mathbf{z} are defined differently. See Section 2.3.3 for the computation of \mathbf{p} .

Model proposed in ISO 6143 The method proposed by ISO 6143 [5] is again a maximum likelihood fit:

$$\max_{\xi, \mathbf{p}} \left(K \exp\left(-\frac{1}{2}(\xi - \mathbf{x})^\top \mathbf{C}_x^{-1}(\xi - \mathbf{x}) - \frac{1}{2}(\mathbf{f}(\xi, \mathbf{p}) - \mathbf{y})^\top \mathbf{C}_y^{-1}(\mathbf{f}(\xi, \mathbf{p}) - \mathbf{y})\right) \right) \quad (2.23)$$

with $K = (2\pi)^{-n} (\det(\mathbf{C}_x) \det(\mathbf{C}_y))^{-1/2}$ constant. Comparing this formula to Equation (2.21), one can see that (2.21) is almost a linearized form of (2.23). But in (2.21) the covariances between calibration standards lead to additional terms that are ignored by (2.23). On the other hand, our model (2.21) uses only linear approximations for the $\mathbf{f}(\xi, \mathbf{p})$ -terms. Both forms have their advantages: the linearized version (2.21) can handle covariances and can be used to compute the probability density function for the P-FIT, whereas the non linearized variant (2.23) is too complicated to proceed on that path. The disadvantage is that one is forced to linearize the calibration function around the calibration standards. The ISO 6143 model is more accurate in cases where the calibration function shows strong non linear behavior in the neighborhood of the calibration standards, this is at the prize of ignoring covariances completely.

In our opinion the benefit of considering covariances is more impor-

tant since in practical cases one strives for almost linear measurement functions. Moreover, C_x and C_y cannot be assumed to be diagonal as pretended by (2.23) (see [5, p. 22]). However, using this assumption anyhow Equation (2.23) leads to the minimization problem

$$\min_{\xi, \mathbf{p}} \left(\sum_{i=1}^n \left(\frac{(\xi_i - x_i)^2}{\alpha_i^2} + \frac{(f(\xi_i, \mathbf{p}) - y_i)^2}{\beta_i^2} \right) \right). \quad (2.24)$$

This is the minimization problem described in [5, p. 19]. It belongs to the class of *weighted orthogonal distance regression* (Weighted ODR) problems that are thoroughly studied in [11] and [43].

2.3.3 Computations

In this section we show how to solve the maximization problems (2.21) and (2.22). For quantifying the measurement uncertainty when using the calibrated tool we need the covariances of the parameters \mathbf{p} , i.e. the covariance matrix C_p . Its use will be demonstrated in Section 4.4. We explain the computation of the covariance matrix C_p at the end of each subsection.

Notation We use the following symbols, notations and abbreviations for readability: $M^{(i)}$ denotes the i -th column vector of matrix M . Additionally, ∂_k abbreviates $\frac{\partial}{\partial p_k}$, i.e. the element-wise derivative with respect to p_k .

$$\begin{aligned} J_{\mathbf{p}}^{\top} &= \left(\frac{\partial}{\partial \mathbf{p}} f^{(1)}, \dots, \frac{\partial}{\partial \mathbf{p}} f^{(n)} \right) \\ J_{x\mathbf{p}}^{\top} &= \left(\frac{\partial}{\partial \mathbf{p}} f_x^{(1)}, \dots, \frac{\partial}{\partial \mathbf{p}} f_x^{(n)} \right) \\ F_{x\mathbf{p}_k} &= \partial_k F_x \end{aligned}$$

Computing ξ and \mathbf{p} Consider the maximization problem (2.21). This is equivalent (after taking the logarithm) to the following minimization problem:

$$\min_{\xi, \mathbf{p}} \left(n \log(2\pi) + \frac{1}{2} \log(\det(C)) + \frac{1}{2} \mathbf{z}^{\top} C^{-1} \mathbf{z} \right).$$

Constant terms can now be dropped. Note that $\det(C)$ is also constant. This can be seen by applying a *Cholesky factorization* to C . Note: There exists a Cholesky factorization for any symmetric positive semi-definite matrix

which, by definition, is the case for all covariance matrices.

$$C = R^\top R = \begin{pmatrix} R_x^\top & \mathbf{0} \\ F_x R_y^\top & R_y^\top \end{pmatrix} \begin{pmatrix} R_x & R_x F_x \\ \mathbf{0} & R_y \end{pmatrix},$$

where $C_x = R_x^\top R_x$ and $C_y = R_y^\top R_y$. Note the fact that $\det(C) = \det(R_x^\top) \det(R_x) \det(R_y^\top) \det(R_y) = \det(R_x^\top R_x) \det(R_y^\top R_y) = \det(C_x) \det(C_y)$ and therefore $\det(C)$ is independent of ξ and p and thus constant. Lucky us! We are left with the following minimization problem

$$\min_{\xi, p} \left(z^\top C^{-1} z \right).$$

Solving this is equivalent to solving a *nonlinear least squares problem*:

$$\begin{aligned} \min_{\xi, p} \left(z^\top C^{-1} z \right) &= \min_{\xi, p} \left(z^\top R^{-1} R^{-\top} z \right) \\ &= \min_{\xi, p} \left(r^\top r \right) \quad \text{with } r = R^{-\top} z \\ &= \min_{\xi, p} \|r\|_2^2 = \min_{\xi, p} \mathcal{S}(\xi, p). \end{aligned} \quad (2.25)$$

There are different standard solvers for such problems, e.g. the *Gauss-Newton* or the further developed *Levenberg-Marquardt* method, other trust region methods, the *Newton* method, the *spectral decomposition* method, etc. [10, 24, 25, 36, 42]. Section 2.3.5 describes the development of the algorithm used in the MUSAC-system to some detail. The following condition holds for the computed minimum

$$J^\top r = \mathbf{0}, \quad (2.26)$$

where J denotes the *Jacobian* of r with respect to ξ and p . Its structure can be seen in Equation (2.28). The statement above can be verified as follows: we define $w = (\xi, p)^\top$. Linearizing the function $r(w + \Delta w) \approx r(w) + J\Delta w$ we get $\mathcal{S}(w + \Delta w) \approx \mathcal{S}(w) + 2\Delta w^\top J^\top r + \mathcal{O}(\|\Delta w\|_2^2)$. Thus Equation (2.26) obviously holds at a minimum.

The Computation of C_p For a good estimation of the measurement error we need the covariance matrix C_p of p . We therefore rewrite Equation (2.26) as

$$h(w, y) = J^\top r = \mathbf{0}. \quad (2.27)$$

The function h is an implicit description of the maximum likelihood criterion. Thus, w is determined by y only (and of course x which is constant for a given set of calibration standards). Since y is a random variable with a known covariance matrix, we are now able to compute the covariance matrix of w and thus of ξ and p .

We define $\bar{w} = (\bar{\xi}, \bar{p}_{\text{est}})^\top = E[w]$ and $\bar{y} = E[y]$ to be the means of the respective random variables. Note that the mean of several estimated parameter vectors \bar{p}_{est} is not necessarily the true and unknown parameter vector. Then

$$w = \bar{w} + \Delta w \quad \text{and} \quad y = \bar{y} + \Delta y.$$

By definition, we may write

$$h(\bar{w}, \bar{y}) = \mathbf{0}$$

and after linearizing h and inserting the equations above we get

$$\begin{aligned} \mathbf{0} &\approx \frac{\partial}{\partial w} h(w, y) \Delta w + \frac{\partial}{\partial y} h(w, y) \Delta y \\ &= H_w \Delta w + H_y \Delta y \end{aligned}$$

Solving for Δw and rearranging the expressions we get

$$\begin{aligned} \Delta w &\approx -H_w^{-1} H_y \Delta y \\ \implies \text{Cov}(w, w) &= E[\Delta w \Delta w^\top] \\ &\approx E[H_w^{-1} H_y \Delta y \Delta y^\top H_y^\top H_w^{-\top}] \\ &= H_w^{-1} H_y (F_x C_x F_x + C_y) H_y^\top H_w^{-\top}. \end{aligned}$$

Note that H_w is the Hessian of S and that its inverse exists if H_w is positive definite which implies the existence of a unique minimum. The covariance matrix C_p corresponds to the lower right $m \times m$ block of $\text{Cov}(w, w)$.

Finally we must show how to compute the matrices H_w and H_y . We start by expanding the expression (2.27) using the variables introduced throughout Section 2.3.2 and 2.3.3. On the one hand, we have

$$\begin{aligned} H_y &= \frac{\partial}{\partial y} h(w, y) \\ &= \left(\frac{\partial}{\partial y} J^\top \right) r + J^\top \left(\frac{\partial}{\partial y} r \right) \end{aligned}$$

$$\begin{aligned}
\mathbf{H}_y &= \mathbf{0} \cdot \mathbf{r} + \mathbf{J}^\top \mathbf{R}^{-\top} \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{R}_x^{-\top} & \mathbf{0} \\ -\mathbf{R}_y^{-\top} \mathbf{F}_x & -\mathbf{R}_y^{-\top} (\text{diag}(\xi - \mathbf{x}) \mathbf{J}_{xp} - \mathbf{J}_p) \end{pmatrix}^\top \mathbf{R}^{-\top} \begin{pmatrix} \mathbf{0} \\ \mathbf{I} \end{pmatrix} \quad (2.28) \\
&= \begin{pmatrix} -\mathbf{F}_x \mathbf{C}_y^{-1} \\ -(\mathbf{J}_{xp}^\top \text{diag}(\xi - \mathbf{x}) - \mathbf{J}_p^\top) \mathbf{C}_y^{-1} \end{pmatrix}.
\end{aligned}$$

On the other hand, for the i -th column $\mathbf{H}_w^{(i)}$, for $1 \leq i \leq n$, i.e. the derivatives with respect to ξ_i , we have

$$\begin{aligned}
\mathbf{H}_w^{(i)} &= \begin{pmatrix} \frac{\partial}{\partial \xi_i} \mathbf{J}^\top \end{pmatrix} \mathbf{r} + \mathbf{J}^\top \begin{pmatrix} \frac{\partial}{\partial \xi_i} \mathbf{r} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & -\mathbf{J}_{xp}^\top \text{diag}(\mathbf{e}_i) \mathbf{R}_y^{-1} \end{pmatrix} \mathbf{r} + \mathbf{J}^\top \mathbf{R}^{-\top} \begin{pmatrix} \mathbf{e}_i \\ \mathbf{0} \end{pmatrix},
\end{aligned}$$

and for $1 \leq k \leq m$, i.e. the derivatives with respect to p_k , we have

$$\begin{aligned}
\mathbf{H}_w^{(n+k)} &= \begin{pmatrix} \partial_k \mathbf{J}^\top \end{pmatrix} \mathbf{r} + \mathbf{J}^\top (\partial_k \mathbf{r}) \\
&= \begin{pmatrix} \mathbf{0} & -\mathbf{F}_{xp_i} \mathbf{R}_y^{-1} \\ \mathbf{0} & -(\mathbf{J}_{xpp_i}^\top \text{diag}(\xi - \mathbf{x}) - \mathbf{J}_{pp_i}^\top) \mathbf{R}_y^{-1} \end{pmatrix} \mathbf{r} \\
&\quad + \mathbf{J}^\top \left(\begin{pmatrix} \mathbf{0} & \mathbf{0} \\ -\mathbf{R}_y^{-\top} \mathbf{F}_{xp_i} & \mathbf{0} \end{pmatrix} \mathbf{z} + \mathbf{R}^{-\top} \begin{pmatrix} \mathbf{0} \\ -\mathbf{J}_p^{(i)} \end{pmatrix} \right),
\end{aligned}$$

where \mathbf{e}_i is the i -th unit vector.

Computing p only Consider the maximization problem (2.22). This is equivalent to the following minimization problem (after taking logarithms again):

$$\min_p \left(\frac{1}{2} n \log(2\pi) + \frac{1}{2} \log(\det(\mathbf{C})) + \frac{1}{2} \mathbf{z}^\top \mathbf{C}^{-1} \mathbf{z} \right).$$

Note that the matrix C used here is different from the one in Section 2.3.3. R is again the upper triangular matrix resulting from the Cholesky factorization of the new C . It is easy to see that $\det(C)$ is no longer constant in p and cannot be neglected. The only constant term negligible is $\frac{1}{2}n \log(2\pi)$. Thus, we are not so lucky anymore! We are now left with the following general minimization problem

$$\min_p \left(\log(\det(C)) + z^T C^{-1} z \right).$$

Solving (2.29) is no longer equivalent to solving a nonlinear least squares problem, as can be seen by the following transformation

$$\begin{aligned} & \min_p \left(\log(\det(C)) + z^T C^{-1} z \right) \\ &= \min_p \left(\log(\det(R^T R)) + z^T R^{-1} R^{-T} z \right) \\ &= \min_p \left(2 \log(\det(R)) + r^T r \right) = \min_p \mathcal{S}(p) \end{aligned} \quad (2.29)$$

Nonetheless the standard techniques mentioned in Section 2.3.3 can be applied to this problem. However, solving Problem (2.29) requires more effort, as the construction of the Jacobian is more difficult (see Section 2.3.3). At the minimum the following condition holds:

$$q + J^T r = 0$$

Here $q = \frac{\partial}{\partial p} \log(\det(R))$ and J denotes the Jacobian of r with respect to p . Again, as in Equation (2.26), this means that the first derivative with respect to p is zero. We show in the following how to compute these entities.

Computing of C_p In the first step we derive an expression for

$$\begin{aligned} q &= \frac{\partial}{\partial p} \log(\det(R)) \\ &= \frac{\partial}{\partial p} \log \left(\prod_{i=1}^n R_{i,i} \right) \\ &= \frac{\partial}{\partial p} \sum_{i=1}^n \log(R_{i,i}). \end{aligned}$$

For the k -th entry of \mathbf{q} , we get

$$\begin{aligned} \mathbf{q}_k &= \partial_k \sum_{i=1}^n \log(\mathbf{R}_{i,i}) \\ &= \sum_{i=1}^n \frac{\partial_k \mathbf{R}_{i,i}}{\mathbf{R}_{i,i}}. \end{aligned}$$

What remains is to show how $\partial_k \mathbf{R}_{i,i}$ can be computed. We begin by observing that

$$\begin{aligned} \partial_k \mathbf{C} &= (\partial_k \mathbf{F}_x) \mathbf{C}_x \mathbf{F}_x + \mathbf{F}_x \mathbf{C}_x (\partial_k \mathbf{F}_x) \\ &= \mathbf{F}_x \mathbf{p}_k \mathbf{C}_x \mathbf{F}_x + \mathbf{F}_x \mathbf{C}_x \mathbf{F}_x \mathbf{p}_k. \end{aligned}$$

With $\partial_k \mathbf{C}$ and with the derivation rules applied to $\mathbf{C} = \mathbf{R}^\top \mathbf{R}$ we get

$$\partial_k \mathbf{C} = \left(\partial_k \mathbf{R}^\top \right) \mathbf{R} + \mathbf{R}^\top (\partial_k \mathbf{R}).$$

This can easily be solved for $\partial_k \mathbf{R}$ and its diagonal elements are the wanted $\partial_k \mathbf{R}_{i,i}$. As in Section 2.3.3 we introduce again an implicit function

$$\mathbf{h}(\mathbf{p}, \mathbf{y}) = \mathbf{q} + \mathbf{J}^\top \mathbf{r} = \mathbf{0}.$$

Due to the complicated structure of \mathbf{C} we abstain from giving an explicit expression for \mathbf{J} , rather we show its computation in a column-wise manner. We write the k -th column of \mathbf{J} by $\mathbf{J}^{(k)}$, for $1 \leq k \leq m$. Then we have

$$\begin{aligned} \mathbf{J}^{(k)} &= \partial_k \mathbf{r} \\ &= \left(\partial_k \mathbf{R}^{-\top} \right) \mathbf{r} + \mathbf{R}^{-\top} (\partial_k \mathbf{r}) \\ &= \left(\partial_k \mathbf{R}^{-\top} \right) \mathbf{r} - \mathbf{R}^{-\top} \mathbf{J} \mathbf{p}^{(k)}. \end{aligned}$$

This leaves us with the task of computing $\partial_k \mathbf{R}^{-\top}$. We have $\mathbf{I} = \mathbf{R}^{-1} \mathbf{R}$. Thus, after taking derivatives on both sides:

$$\begin{aligned} \mathbf{0} &= (\partial_k \mathbf{R}^{-1}) \mathbf{R} + \mathbf{R}^{-1} (\partial_k \mathbf{R}) \\ \implies \partial_k \mathbf{R}^{-1} &= -\mathbf{R}^{-1} (\partial_k \mathbf{R}) \mathbf{R}^{-1} \end{aligned}$$

Again, as in Section 2.3.3, we need to compute \mathbf{H}_p and \mathbf{H}_y . Let $\mathbf{H}_y^{(i)}$ denote the i -th column of \mathbf{H}_y . For $1 \leq i \leq n$, we get

$$\begin{aligned} \mathbf{H}_y^{(i)} &= \frac{\partial}{\partial y_i} (\mathbf{q} + \mathbf{J}^\top \mathbf{r}) \\ &= \mathbf{0} + \left(\frac{\partial}{\partial y_i} \mathbf{J}^\top \right) \mathbf{r} + \mathbf{J}^\top \left(\frac{\partial}{\partial y_i} \mathbf{r} \right) \\ &= \left(\frac{\partial}{\partial y_i} (\mathbf{J}^{(1)}, \dots, \mathbf{J}^{(m)})^\top \right) \mathbf{r} + \mathbf{J}^\top \mathbf{R}^{-\top} \mathbf{e}_i \\ &= \left((\partial_1 \mathbf{R}^{-\top}) \mathbf{R}^{-\top} \mathbf{e}_i, \dots, (\partial_m \mathbf{R}^{-\top}) \mathbf{R}^{-\top} \mathbf{e}_i \right)^\top \mathbf{r} + \mathbf{J}^\top \mathbf{R}^{-\top} \mathbf{e}_i. \end{aligned}$$

The derivation for \mathbf{H}_p is exceedingly odd. Nevertheless here it is: First we have

$$\partial_l \mathbf{q}_k = \sum_{i=1}^n \frac{(\partial_l \partial_k \mathbf{R}_{i,i}) \mathbf{R}_{i,i} - (\partial_k \mathbf{R}_{i,i}) (\partial_l \mathbf{R}_{i,i})}{\mathbf{R}_{i,i}^2},$$

$(\partial_l \partial_k \mathbf{R}_{i,i})$ can be found by solving

$$\partial_l \partial_k \mathbf{C} = \left(\partial_l \partial_k \mathbf{R}^\top \right) \mathbf{R} + \left(\partial_l \mathbf{R}^\top \right) (\partial_k \mathbf{R}) + \left(\partial_k \mathbf{R}^\top \right) (\partial_l \mathbf{R}) + \mathbf{R}^\top (\partial_l \partial_k \mathbf{R}).$$

Further we need

$$\begin{aligned} \partial_l \mathbf{J}^\top \mathbf{r} &= (\partial_l \mathbf{J})^\top \mathbf{r} + \mathbf{J}^\top (\partial_l \mathbf{r}) \\ &= \left(\partial_l \mathbf{J}^{(1)}, \dots, \partial_l \mathbf{J}^{(m)} \right)^\top \mathbf{r} + \mathbf{J}^\top (\partial_l \mathbf{r}) \end{aligned}$$

$$\partial_l \mathbf{r} = \left(\partial_l \mathbf{R}^{-\top} \right) \mathbf{z} - \mathbf{R}^{-\top} \mathbf{J}_p^{(l)}$$

$$\partial_l \mathbf{J}^{(k)} = \left(\partial_l \partial_k \mathbf{R}^{-\top} \right) \mathbf{r} + \left(\partial_k \mathbf{R}^{-\top} \right) (\partial_l \mathbf{r}) - \left(\partial_l \mathbf{R}^{-\top} \right) \mathbf{J}_p^{(k)} - \mathbf{R}^{-\top} \left(\partial_l \mathbf{J}_p^{(k)} \right).$$

Putting things together yields

$$\mathbf{H}_p^{(l)} = \partial_l \mathbf{q} + \partial_l \mathbf{J}^\top \mathbf{r},$$

where $\partial_l \mathbf{q}$ is of course the vector with elements $\partial_l \mathbf{q}_{k'}$ for $k = 1, \dots, m$.

2.3.4 Measurement

One goal of calibration is to compute a good set of parameters to quantify the measurement tool's behavior. An equally important goal is to estimate the measurement uncertainty when using the calibrated device.

Assume that the measuring tool is calibrated. Thus we have a set of estimated parameters $\hat{\mathbf{p}}$ together with an estimated covariance matrix $\hat{\mathbf{C}}_{\mathbf{p}} = \text{Cov}(\hat{\mathbf{p}}, \hat{\mathbf{p}})$.

Let $\mathbf{y} \in \mathbb{R}^s$ be a measurement vector whose elements are readings taken from the measurement tool. The readings are subject to error of the form $\mathbf{y} - \bar{\mathbf{y}} \sim \mathcal{N}(\mathbf{0}, \mathbf{E})$, where $\bar{\mathbf{y}}$ is the correct but unknown reading vector and $\mathbf{E} = \text{Cov}(\mathbf{y}, \mathbf{y})$.

For $1 \leq i \leq s$ compute \hat{x}_i by applying $\hat{x}_i = g(y_i, \hat{\mathbf{p}})$ or by solving $y_i = f(\hat{x}_i, \hat{\mathbf{p}})$. For arbitrary functions $f(x, \hat{\mathbf{p}})$, differentiable at least once, we propose the use of a standard solver such as *bisection*, *regula falsi*, *Newton's method*, whichever you prefer.

The computation of the measurement uncertainty is performed as follows. We represent the physical law by $\bar{y}_i = f(\bar{x}_i, \mathbf{p}_{\text{ex}})$, where \bar{x}_i , \mathbf{p}_{ex} and \bar{y}_i are the exact but unknown "true values". The given values are $\bar{x}_i + \Delta x_i = \hat{x}_i$, $\bar{y}_i + \Delta y_i = y_i$ and $\mathbf{p}_{\text{ex}} + \Delta \mathbf{p} = \hat{\mathbf{p}}$. Substituting and linearizing these equations leads to

$$\Delta y_i \approx f_x^{(i)} \Delta x_i + \nabla_{\mathbf{p}} f^{(i)\top} \Delta \mathbf{p},$$

where $f^{(i)} = f(\hat{x}_i, \hat{\mathbf{p}})$ and $f_x^{(i)} = f_x(\hat{x}_i, \hat{\mathbf{p}})$. Extending this approximation to s readings \mathbf{y} and solving for $\Delta \mathbf{x}$, we get

$$\Delta \mathbf{x} \approx \mathbf{F}_x^{-1} \Delta \mathbf{y} - \mathbf{F}_x^{-1} \mathbf{J}_{\mathbf{p}} \Delta \mathbf{p}$$

and furthermore

$$\begin{aligned} \mathbf{C}_x &= \mathbb{E} \left[\Delta \mathbf{x} \Delta \mathbf{x}^\top \right] \\ &\approx \mathbf{F}_x^{-1} \text{Cov}(\mathbf{y}, \mathbf{y}) \mathbf{F}_x^{-\top} \\ &\quad + \mathbf{F}_x^{-1} \mathbf{J}_{\mathbf{p}} \text{Cov}(\hat{\mathbf{p}}, \hat{\mathbf{p}}) \mathbf{J}_{\mathbf{p}}^\top \mathbf{F}_x^{-\top} \\ &\quad - \mathbf{F}_x^{-1} \left(\text{Cov}(\mathbf{y}, \hat{\mathbf{p}}) \mathbf{J}_{\mathbf{p}}^\top + \mathbf{J}_{\mathbf{p}} \text{Cov}(\hat{\mathbf{p}}, \mathbf{y}) \right) \mathbf{F}_x^{-\top}. \end{aligned}$$

There are two things to note. First, $\text{Cov}(\mathbf{y}, \hat{\mathbf{p}})$ can safely be assumed to be $\mathbf{0}$, since the uncertainties of the reading for the s measurements are indepen-

dent of $\hat{\boldsymbol{p}}$ found during calibration. Second, it is obvious that not a single $f_x^{(i)}$ is allowed to be 0. In practice this does not represent a problem, since no-one would use a measurement device within a range where the slope of the calibration function is even close to 0. So we end up with the calculation for

$$\mathbf{C}_x \approx \mathbf{F}_x^{-1} \left(\text{Cov}(\mathbf{y}, \mathbf{y}) + \mathbf{J}_p \text{Cov}(\hat{\boldsymbol{p}}, \hat{\boldsymbol{p}}) \mathbf{J}_p^\top \right) \mathbf{F}_x^{-1}. \quad (2.30)$$

2.3.5 Truncated SVD Solver

In a first approach we applied an ordinary *Gauss-Newton* method to solve the Equations (2.24), (2.25) or (2.29). But we soon realized that the nonlinearities of the problems were strong enough to render the reliability of the plain Gauss-Newton method insufficient. We then applied a *Levenberg-Marquardt* method (see Algorithm 5) which stabilizes the Gauss-Newton method in the sense that it converges more reliably at the price of more iterations [42]. The effect of the λ introduced is to reduce the norm of the correction vector \boldsymbol{d} , such that the norm of the residual \boldsymbol{r} decreased in every step. With increasing λ the norm of \boldsymbol{d} diminishes. If $\lambda = 0$ then Algorithm 5 corresponds to exactly the Gauss-Newton method. There are uncountable proposals for a strategy to adapt λ . None of them (at least the ones we experimented with) was satisfactory in the range of problems we were facing in the MUSAC-project. It turned out that problems stemming from calibration have surprisingly often a badly conditioned Jacobian \mathbf{J} . This can lead to bad numerical results for \boldsymbol{d} , and therefore the whole algorithm is disappointingly unstable, numerically. The Levenberg-Marquardt step solves $\mathbf{J}\boldsymbol{d} \approx -\boldsymbol{r}$ with some constraints on $\|\boldsymbol{d}\|_2$.

Alternatively we can compute an unconstrained step \boldsymbol{d} using the *Singular Value Decomposition* (SVD). Then the following steps are executed (see e.g. [27]):

$$\begin{aligned} \mathbf{J} &= \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top & \mathbf{U} \text{ and } \mathbf{V} \text{ orthogonal, } \boldsymbol{\Sigma} &= \text{diag}(\sigma_1, \dots, \sigma_n) \\ \tilde{\boldsymbol{r}} &= \mathbf{U}^\top \boldsymbol{r} & \text{apply left singular vectors} \\ \forall i : \sigma_i \neq 0 : \tilde{\boldsymbol{d}}_i &= \tilde{r}_i / \sigma_i & \text{compute the solution} \\ \boldsymbol{d} &= \mathbf{V}\tilde{\boldsymbol{d}} & \text{apply right singular vectors} \end{aligned}$$

This procedure uses the computationally expensive SVD, but it discloses useful structural information about the system of equations. The procedure

Algorithm 5 Gauss-Newton Loop with Levenberg-Marquardt Method for Solving Equation (2.25)

function $\langle \xi, p \rangle = \text{LM_solver}(x, y, R_y, R_x)$
 — Find ξ and p minimizing (2.25) or similar —
 $\langle \xi, p, \lambda \rangle \leftarrow \text{Initialize}(x, y)$
do — Gauss-Newton-Loop —
 $z \leftarrow [\xi - x; y - f(\xi, p)]$ — set up z, r and J —
 $r \leftarrow R^{-\top} z$
 $J \leftarrow \begin{pmatrix} R_x^{-\top} & 0 \\ -R_y^{-\top} F_x & -R_y^{-\top} (\text{diag}(\xi - x) J_{xp} - J_p) \end{pmatrix}$
 $d \leftarrow -(J^{\top} J + \lambda I)^{-1} (J^{\top} r)$ — Levenberg-Marquardt step and adaption
 $\lambda \leftarrow \text{Adapt}(\lambda)$ of λ according to some heuristics, to be
 replaced by Algorithm 6 —
 $[\xi; p] \leftarrow [\xi; p] + d$ — update ξ and p —
until $\|Jr\|_2$ has converged

does not impose any constraints on d . But note that the division by the singular values σ_i can lead to ridiculously large entries in \tilde{d} depending on the respective entry in r and the magnitude of σ_i . Even worse, if any σ_i happens to be 0 this procedure would lead to a crash. There is, again, a whole bunch of proposals on how to *truncate* the SVD in this context to ensure a good convergence for Algorithm 5. A good survey on trust region methods is [35]. We experimented with some of them. In the end, we came up with a truncation that we haven't seen anywhere else before. It seems astonishingly well suited for the class of problems stemming from calibration. It is shown in Algorithm 6 on the following page. The key idea of Algorithm 6 is that the norm of the correction vector d is bounded by a fraction λ of the current solution's norm $\|z\|_2$. Note that $\|\tilde{d}\|_2 = \|d\|_2$ due to the orthogonality of V . Therefore after the loop over all i , the number n is equal to $\|d\|_2^2$. Within that loop, we observe whether $\|d\|_2 > \lambda \|z\|_2$. If this is the case then the boolean vector b stores the entries of \tilde{d} responsible for the norm exceeding the trust region. After the loop, if any of b 's elements were marked, the factor α is computed to be the factor by which to reduce the exceeding entries in \tilde{d} , such that $\|\tilde{d}\|_2 \leq \lambda \|z\|_2$ holds.

Algorithm 6 *Truncated SVD Method*

```

function  $\langle d, \lambda \rangle = \text{TSVD}(J, r, z, \lambda)$ 
  — Find correction  $d$  for solution so far  $z$  and adapt  $\lambda$  —
   $\langle U, \Sigma, V \rangle \leftarrow \text{svd}(J)$  — apply SVD on  $J$  —
   $\tilde{r} \leftarrow U^T r$  — apply left singular vectors —
   $n \leftarrow 0, b \leftarrow 0$ 
  for all  $i$  do
    if  $\sigma_i \neq 0$  then
       $\tilde{d}_i \leftarrow \tilde{r}_i / \sigma_i$ 
       $n \leftarrow n + \tilde{d}_i^2$  — update the norm of  $d$  —
      if  $n > \lambda^2 \|z\|_2^2$  then
         $b_i \leftarrow 1$ 
      end if
    else
       $\tilde{d}_i \leftarrow 0$ 
    end if
  end do
  if any( $b = 1$ ) then
     $\alpha \leftarrow \frac{\lambda \|z\|_2 - \|(1 - b_i)\tilde{d}_i\|_2}{\|b_i\tilde{d}_i\|_2}$  — elementwise multiplications —
     $\tilde{d}_i \leftarrow (1 - (1 - \alpha)b_i)\tilde{d}_i$  — scale down the entries of  $\tilde{d}$  that
    — exceed the allowed norm of  $d$  —
  end if
   $d \leftarrow V\tilde{d}$  — apply right singular vectors —
   $\lambda \leftarrow \beta\lambda$  — our heuristic to adapt  $\lambda$ :  $\beta < 1$  —

```

Many variants of TSVD propose to ignore singular values below some heuristically adapted threshold, few propose to consider the quotient r_i/σ_i and ignore it if some heuristic criteria are met. We do not ignore any of the information gained from the SVD, but we scale down the seemingly over exaggerated contributions of either r_i and/or σ_i .

Often people propose to reduce the length of the correction so that the norm really diminishes in every iteration of the Gauss-Newton-Loop. We do not require this. Since $\tilde{\mathbf{d}}$ computed by Algorithm 6 is not even guaranteed to be a descending direction, despite the fact that it is certainly an approximation of an descending direction. By reducing λ monotonically, we ensure that the correction eventually gets so small, that the entire $\tilde{\mathbf{d}}$ gets scaled, instead of only parts of it. Now $\alpha\mathbf{V}\tilde{\mathbf{d}}$ is indeed a true descending direction, and therefore (again when λ and thus α is small enough) the minimizing Gauss-Newton-Loop is really progressing towards the next local minimum (at least).

In the MUSAC-system we initialize λ to 0.75 and the reduction factor β to 0.991. This yields 4044 iterations until $\lambda < 10^{-16}$. If this many iterations are used without $\|\mathbf{J}\mathbf{r}\|_2$ ever getting sufficiently close to 0 then a warning about non-convergence is issued and the Gauss-Newton-Loop is aborted. It is difficult to compare such heuristics thoroughly. Because any comparison depends strongly on the input given to the solver. Therefore the following comparison of the TSVD-solver with the Levenberg-Marquardt solver (with a fixed λ) is far from complete. It is a data set of an arbitrary calibration, as it may appear in the HPLC:

$$\begin{aligned} \mathbf{x} &= (100, 10, 1, 0.1)^\top \\ \mathbf{y} &= (10000, 1000, 100, 10)^\top \\ \mathbf{C}_x &= \begin{pmatrix} 133.333 & 13.3333 & 1.33333 & 0.133333 \\ 13.3333 & 1.33348 & 0.133348 & 0.0133348 \\ 1.33333 & 0.133348 & 0.0133362 & 0.00133362 \\ 0.133333 & 0.0133348 & 0.00133362 & 0.000133376 \end{pmatrix} \quad (2.31) \\ \mathbf{C}_y &= \text{diag}(10000, 100, 1, 0.0001) \end{aligned}$$

This calibration input data was used to fit a line $f(x) = p_0 + p_1x$ with the *Classical Model* which is used in the Monte Carlo Simulation, and to *Berkson's Model*, the XiP-FIT employed in the Iso-Evaluation. Table 2 shows the average of iterations used by each of the solvers to reach a solution within the tolerance of 10^{-8} . The starting parameter vector \mathbf{p}_0 was varied from

$\mathbf{p} = (p_0, p_1) \in \{-100, -10, 0, 1, 10, 100\}^2$. The correct \mathbf{p} depends on the model, but in both cases it is very close to $\mathbf{p} = (0.0, 100.0)^\top$. Table 2 shows that in terms of iterations until convergence the two solver are equally fast. But for Berkson's Model which is intrinsically worse conditioned, the TSVD solver is faster, tremendously so. Further experiments show also that the TSVD solver is much more robust against bad starting points. These are two advantages of the proposed TSVD solver that outweigh the computationally more expensive singular value decomposition by far.

Table 2 Average Iterations for TSVD-Solver and Levenberg-Marquardt Solver on the Input Data (2.31) at a Tolerance of 10^{-8}

	TSVD	Levenberg-Marquardt
Classical Model	$\mu = 7.0, \sigma = 1.4$	$\mu = 7.0, \sigma = 1.4$
Berkson's Model	$\mu = 4.9, \sigma = 1.6$	$\mu = 3678, \sigma = 1074$

2.3.6 Numerical Comparisons and Simulations

In the following section we compare four different calibration functions. We fit a linear, an exponential, a power and a sigmoidal function. In these experiments we compare the data to a "God given", known and exact function. The calibration data was generated randomly according to Figure 2 on page 32 using the true parameters \mathbf{p}_{ex} . The covariances in x and y are assumed known as well. The variance in y at measurement time was chosen quite arbitrarily.

Simply put, we are trying to recover the true function $f(x, \mathbf{p}_{\text{ex}})$. One would expect that the chances for recovering the true function are bad using only a handful of calibration points. But the parameters recovered or fitted by the different methods are so close that, in order to make the differences visible at all, we had to apply big variances for the random errors and we selected the calibration points at unfavorable x -positions for the entire calibration process.

In each case the figure on the left hand side shows the fitted functions using XIP-FIT (thick line), P-FIT (thin line), the classical least squares fit (dash-dotted line) and the exact line as a fine dotted line. The two dashed lines

build an approximated 95%–confidence region for the XIP-FIT. The confidence region is computed by constructing 1000 random instances of the fitted parameters within their 95%–confidence interval. The dashed line is then the hull of all curves that are produced by each of these 1000 instances.

The figure on the right shows the simulated measurement uncertainty for the XIP-FIT and the P-FIT (dashed lines) compared to the estimated uncertainties (solid lines) according to Equation (2.30). The simulated uncertainty is computed as follows: Each of $\text{Cals} = 200$ random instances of $f^{-1}(y, \mathbf{p})$ are applied to $\text{Sims} = 20$ random instances of y -values. This yields $\text{Cals} \times \text{Sims} = 4000$ measured x -values whose variance is plotted at the mean position of those x -values. This process is done for 20 y -positions within the range that should according to \mathbf{p}_{ex} lay in the x -range of the calibration standards

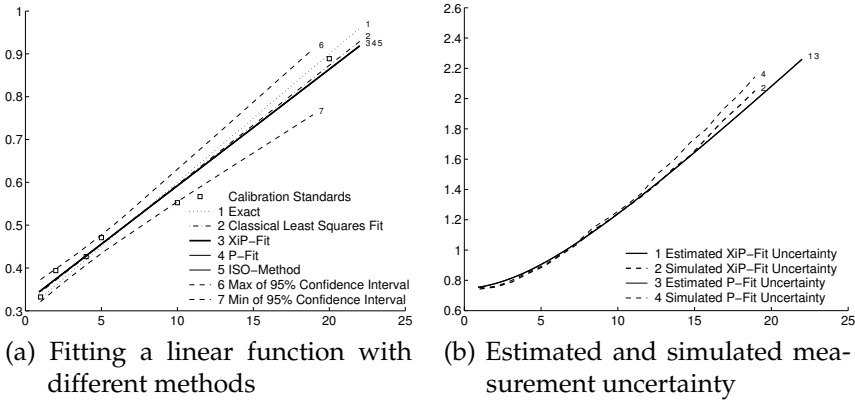
Linear Fit The calibration function is $f(x, \mathbf{p}) = p_1 + p_2x$. The calibration data for Figure 3 is

x	1.0000	2.0000	4.0000	5.0000	10.0000	20.0000
y	0.332295	0.393916	0.426571	0.470667	0.552436	0.888801

The covariance in x is given by the following matrix:

$$C_x = \begin{pmatrix} 0.0012 & 0.0016 & 0.0016 & & & & \\ 0.0016 & 0.0032 & 0.0032 & & & & \\ 0.0016 & 0.0032 & 0.0064 & & & & \\ & & & 0.0100 & & & \\ & & & & 0.0400 & & \\ & & & & & 0.1600 & \end{pmatrix}$$

In other words, their relative standard deviation is around 2% and the first three x -values are correlated. The covariance in y is $C_y = 10^{-4} \text{diag}(2.7225, 3.24, 4.41, 5.0625, 9, 20.25)$. In other words, they are not correlated and their relative standard deviation is in the order of 5%. Finally the fitted parameters \mathbf{p} are as follows:

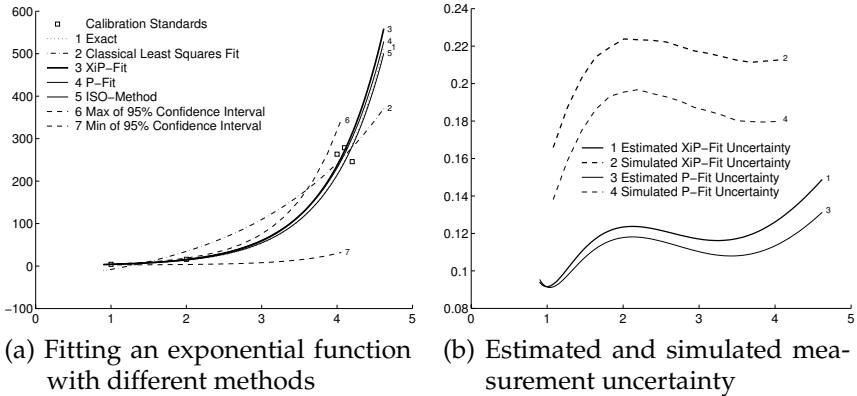
Figure 3 Fitting a linear function

	p_1	p_2
Exact	0.3	0.03
Classical LS	0.31558	0.027886
XIP-FIT	0.31971	0.027226
P-FIT	0.31971	0.027204
ISO-Method	0.31972	0.027226

The ISO-Method as sketched on page 35 is implemented in B-Least, a software available at BAM.

Remarks The matching between estimated and simulated measurement uncertainty is quite nice. The difference between the compared methods is not very big, in fact for linear functions and if the x -values are not correlated then the XIP-FIT and the ISO-Method can be considered equivalent. The difference in the first parameter is entirely due to the correlation of the x -values. The variation of the second parameter, the slope p_2 between the P-FIT and XIP-FIT is due to the different model assumptions. Note the big uncertainties involved in this example: the y -values for the calibration standards are all assumed to have a standard deviation of 5%. For a measurement of $y = 0.4$, $u_y = 5\%$ this yields an estimated $x = 2.94917$ with $u_x = 0.809 = 27.5\%$! Example 18 on page 89 shows the M input using the data of this example.

Figure 4 Fitting an Exponential Function



Exponential Fit The calibration function is $f(x, p) = p_1 + p_2 \exp(p_3 x)$. The data for Figure 4 is

x	1.0000	2.0000	4.0000	4.1000	4.2000
y	4.1682	15.9362	263.3923	278.6612	245.8233

The covariance in x is given by the following matrix:

$$C_x = \begin{pmatrix} 0.0075 & 0.0100 & 0.0100 & & \\ 0.0100 & 0.0200 & 0.0200 & & \\ 0.0100 & 0.0200 & 0.0400 & & \\ & & & 0.0420 & \\ & & & & 0.0441 \end{pmatrix}$$

The covariance in y is $C_y = \text{diag}(0.0280, 0.4393, 117.1169, 154.9422, 204.9871)$. Finally the fitted parameters p are as follows:

	p_1	p_2	p_3
Exact	0.1	0.8	1.4
Classical LS	-62.9821	31.2164	0.5699
XiP-FIT	0.39654	0.91511	1.3883
P-FIT	0.47711	0.88251	1.3724
ISO-Method	-0.060929	1.1143	1.3338

Remarks The uncertainties in the calibration and the measurement data are 5%. This is quite considerable. The errors induced by linearization now become obvious. We presume that the simulated uncertainties are therefore bigger but qualitatively they are still acceptable. Note also that the classical least squares fit is way out of bounds. The ISO-Method, despite the different parameter numbers, crowds the picture in between the XIP-FIT and the P-FIT. There was a warning of the B-Least software about a division by zero but the output still seems reasonable. The picture on the right was produced with considerably more simulations as in the linear case above. (Cals = 1000, Sims = 10).

Power Function Fit The calibration function is $f(x, \mathbf{p}) = p_1 x^{p_2}$. Some of the data for Figure 5 is taken from NIST¹.

x	1.3000	1.4700	1.4900	1.5700	1.6000	1.6800
y	2.3093	2.6972	3.9351	3.8846	4.9255	6.2209

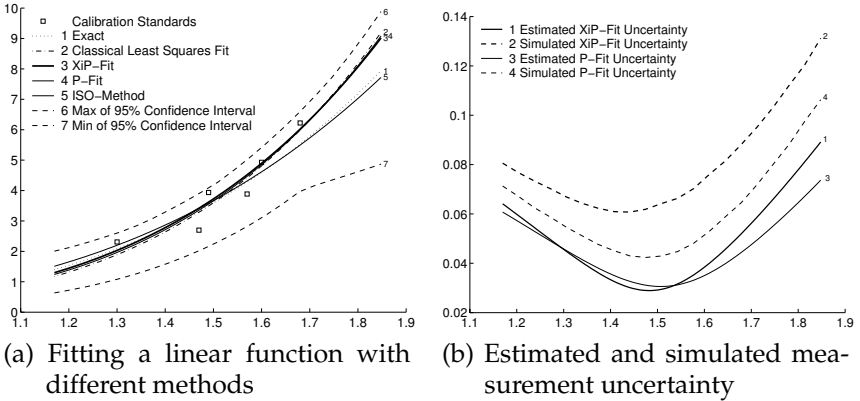
The covariance in x is $C_x = \text{diag}(0.0027, 0.0035, 0.0036, 0.0039, 0.0041, 0.0045)$. The covariance in y is $C_y = \text{diag}(0.0070, 0.0177, 0.0196, 0.0292, 0.0338, 0.0489)$. Finally the fitted parameters \mathbf{p} are as follows:

	p_1	p_2
Exact	0.7700	3.8000
Classical LS	0.5801	4.5005
XIP-FIT	0.6676	4.2387
P-FIT	0.8618	3.5703
ISO-Method	0.6264	4.3545

Remarks The uncertainties in the calibration and the measurement data are 4%. Again the linearization induces a shift between the estimated and the simulated measurement uncertainty (Cals = 1000, Sims = 10). Note that the P-FIT is visually better in the sense that it is closer to the exact function. This is the case for most of the experiments we have studied. The computation for the ISO-method is done by our own MATLAB-implementation because the B-Least software does not offer this type

¹<http://www.itl.nist.gov/div898/strd/index.html> → Nonlinear Regression → Dan Wood

Figure 5 Fitting a Power Function



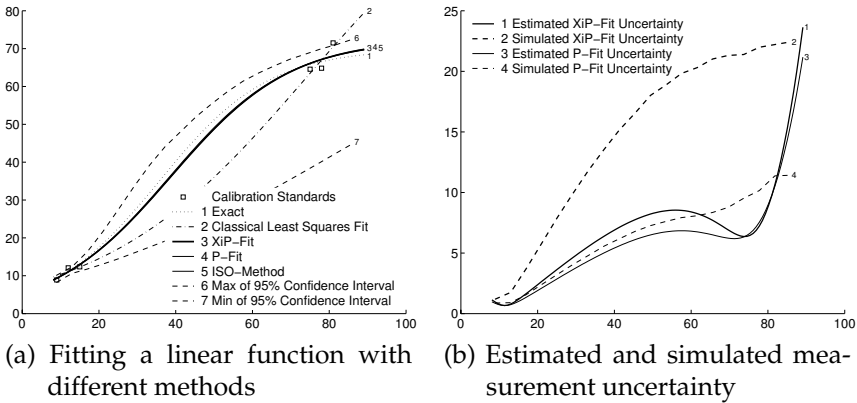
of function to be fitted. Despite numerous trials we couldn't find an input data set which would have shown more obvious differences between the methods. Our conclusion is that for fitting a power function the XIP-FIT and the ISO-method are practically equal, whereas the P-FIT seems to be slightly better recovering the true function.

Sigmoid Function Fit The calibration function is $f(x, \mathbf{p}) = p_1 / (1 + \exp(p_2 - p_3 x))$. Some of the data for Figure 6 is taken from NIST². They assign a "higher level of difficulty" to this problem.

x	9	12	15	75	78	81
y	8.8916	12.0857	12.3492	64.5600	64.8196	71.5168

The covariance in x is $C_x = \text{diag}(0.0729, 0.1296, 0.2025, 5.0625, 5.4756, 5.9049)$. The covariance in y is $C_y = \text{diag}(0.0786, 0.1126, 0.1592, 3.8960, 3.9862, 4.0616)$. The fitted parameters \mathbf{p} are as follows:

²<http://www.itl.nist.gov/div898/strd/index.html> → Nonlinear Regression → Rat42

Figure 6 Fitting a Sigmoid Function

	p_1	p_2	p_3
Exact	70.0000	2.5000	0.0700
Classical LS	124.6656	2.7761	0.0375
XiP-FIT	72.5830	2.4889	0.0641
P-FIT	72.2162	2.4931	0.0650
ISO-Method	72.6064	2.4900	0.0641

Remarks The uncertainties in the calibration and the measurement data are 3%. As opposed to the other examples, the estimated and the simulated measurement uncertainties for the P-FIT are very close. We attribute this to the fact that the uncertainties in the calibration data are the same and relatively small. On the other hand the XiP-FIT seems to yield uncertainty estimates which are too optimistic. For this plot we used $Ca_{ls} = 500$ and $Sims = 10$. Note how the measurement uncertainty seems to explode as soon as the calibration function's slope approaches zero.

2.3.7 Final remarks on Regression

Taking the covariance of calibration data into account is important for cases with very strong correlation or with relatively large uncertainties. However,

there is no sound reason to ignore it, since it does not seriously affect the runtime behavior.

The P-FIT leads to slightly higher measurement uncertainties, but it recovers the “exact” parameters generally better — at least in our experiments. This is explicable through its property of “integrating over all possible true ξ -values”. On the other hand it is computationally more expensive. The XIP-FIT yields the same results for linear regressions as the ISO-Method as long as no correlations among the x -values are involved. When the function to be fitted is not linear, then the XIP-FIT and the ISO method differ due to the fact that the XIP-FIT linearizes the function taking covariances into account and the ISO method does not linearize the function but ignores covariances. Classical least squares fitting is not generally suitable for calibration tasks. There may be exceptions to this rule in peculiar situations where one knows and/or is well aware about negligible variances.