

# Statistics of Parameter Estimates: A Concrete Example\*

Oscar Aguilar<sup>†</sup>  
Moritz Allmaras<sup>‡</sup>  
Wolfgang Bangerth<sup>§</sup>  
Luis Tenorio<sup>¶</sup>

**Abstract.** Most mathematical models include parameters that need to be determined from measurements. The estimated values of these parameters and their uncertainties depend on assumptions made about noise levels, models, or prior knowledge. But what can we say about the validity of such estimates, and the influence of these assumptions? This paper is concerned with methods to address these questions, and for didactic purposes it is written in the context of a concrete nonlinear parameter estimation problem. We will use the results of a physical experiment conducted by Allmaras et al. at Texas A&M University [M. Allmaras et al., *SIAM Rev.*, 55 (2013), pp. 149–167] to illustrate the importance of validation procedures for statistical parameter estimation. We describe statistical methods and data analysis tools to check the choices of likelihood and prior distributions, and provide examples of how to compare Bayesian results with those obtained by non-Bayesian methods based on different types of assumptions. We explain how different statistical methods can be used in complementary ways to improve the understanding of parameter estimates and their uncertainties.

**Key words.** parameter estimation, Bayesian inference, frequentist inference, model validation, data analysis, residual analysis, maximum likelihood, nonlinear regression, surrogate models

**AMS subject classifications.** 34A55, 62F15, 62P35, 97M50

**DOI.** 10.1137/130929230

**I. Introduction.** Mathematical models of physical phenomena usually depend on parameters that need to be estimated from measurement data. Such estimation can be done using a variety of methods, but since no single procedure works best for every problem, we try to select a method such that (i) it takes advantage of all relevant available information; (ii) it is computationally feasible; (iii) it makes no unreasonable assumptions; and (iv) it is consistent with the data and physical models. However, we are sometimes too eager to focus on (i) and (ii) and neglect (iii) and (iv), which are important because the validity of the results depends on satisfying the assumptions of the statistical model. In particular, reported uncertainties are meaningful to the

\*Received by the editors July 15, 2013; accepted for publication (in revised form) April 10, 2014; published electronically February 5, 2015.

<http://www.siam.org/journals/sirev/57-1/92923.html>

<sup>†</sup>Department of Statistics, Iowa State University, Ames, IA 50011 (oaguilar@iastate.edu).

<sup>‡</sup>Siemens AG, Otto-Hahn-Ring 6, 81739 Munich, Germany (moritz.allmaras@gmail.com).

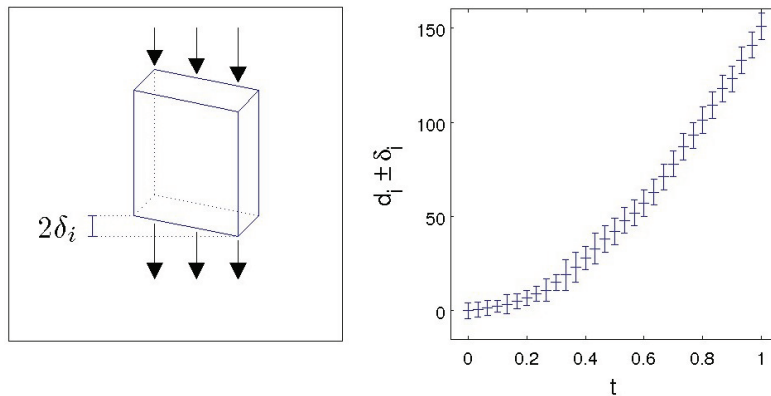
<sup>§</sup>Department of Mathematics, Texas A&M University, College Station, TX 77843 (bangerth@math.tamu.edu). The work of this author was partially supported by award KUS-C1-016-04, made by King Abdullah University of Science and Technology (KAUST).

<sup>¶</sup>Department of Applied Mathematics and Statistics, Colorado School of Mines, CO 80401 (ltenorio@mines.edu). The work of this author was partially supported by NSF grant DMS-0914987.

extent that the assumptions are approximately satisfied. In this paper we describe different ways to check assumptions and explore the validity of parameter estimates and their uncertainties. This is done in the context of a concrete nonlinear parameter estimation problem. We use the results of the experiment published by Allmaras et al. in this journal in 2013 [1]. Thus, this article is a continuation of that tutorial, which we will henceforth refer to as “A13.”

The experiment discussed in A13 consisted of recording the free fall of a box dropped from a known height. The fall was recorded by a video camera at a rate of thirty frames per second with the aim of determining the gravitational acceleration and air friction coefficient of the box. The distance traveled by the object was obtained by analyzing 31 frames of this video.  $t_0$  is the time in seconds when the object starts to travel and was initially chosen to be the time corresponding to the first of the 31 frames, with  $t_i$  the times corresponding to the other frames. As was shown in A13, the time,  $t_0$ , when the object actually started to travel was impossible to determine accurately because any motion in the first frames is small compared with the resolution of the camera, and the second part of A13 therefore treats it as uncertain as well.

From the frames of the video showing the body falling past a tape measure, the authors determined the distance  $d_i$  the body has fallen at time  $t_i$ . But, of course, measuring  $d_i$  is subject to uncertainties. Even more than motion blur, the most important source of error in these data was the fact that the box rotated during its fall and the consequential uncertainty in how to define the location of the object. A13 defined the error,  $\delta_i$ , to be half the distance between the lower and upper corners of the bottom of the box,<sup>1</sup> as depicted in Figure 1.1 (see also the photographs in Figure 3.2 of A13).



**Fig. 1.1** *Left: Depiction of the falling box and the errors  $\delta_i$ . See Figure 3.2 of A13 for an actual picture of the falling box in the experiments described there. Right: Plot of  $d_i \pm \delta_i$ .*

These values were used in A13 to define the support of the probability density of the errors associated with  $d_i$ . The right panel in Figure 1.1 shows the distance  $d_i$  (in inches) as a function of  $t_i$  (s) with error bars defined as  $\pm\delta_i$ .

<sup>1</sup>This might not be the best definition of the measurement error, given that the box can rotate around all three axes, but it is one that could be extracted from the individual frames of the video. As we will discuss below, their choice was overly conservative, and so the details of their definition of measurement error might not have mattered very much.

As explained in A13, Newton's equations of motion lead to an explicit expression of the distance  $z(t)$  traveled by the object at time  $t$ :

$$(1.1) \quad z(t) = (1/C) \log \cosh \left[ \sqrt{gC}(t - t_0) \right],$$

where  $g$  is the acceleration due to gravity and  $C$  is the specific coefficient of air resistance for this particular box. Although both parameters are unknown, our main interest is in identifying  $g$ , for which we have the reliable estimate

$$(1.2) \quad g \approx g_{\text{ref}} = 9.7935 \text{ (m/s}^2\text{)}$$

for the location where the experiment took place (College Station, TX; see A13). Here, as in A13, we are therefore primarily concerned with estimating  $g$  and comparing it against  $g_{\text{ref}}$  for validation. As in every experiment, we do not of course use this knowledge of a reference value in the estimation procedures. Since A13 determined that knowledge of the initial time  $t_0$  was subject to a small error, it considered the case where  $t_0$  is unknown. Hence, the vector of unknown parameters we use will be either  $\mathbf{m} = (g, C)$  or  $\mathbf{m} = (g, C, t_0)$ .

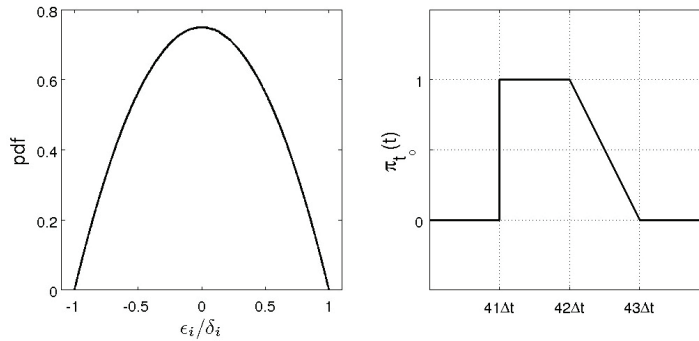
To summarize, the data,  $d_i$ , consist of a noisy recording of  $z(t)$  at 31 different times  $t_i$ :  $d_i = z(t_i) + \varepsilon_i$ , where  $\varepsilon_i$  is modeled as random noise. We write the data vector as

$$(1.3) \quad \mathbf{d} = \mathbf{z}(\mathbf{m}) + \boldsymbol{\varepsilon}, \quad z_i(\mathbf{m}) = (1/C) \log \cosh \left[ \sqrt{gC}(t_i - t_0) \right].$$

The objective of the inverse problem is to use the data  $\mathbf{d}$  to estimate the unknown parameter vector  $\mathbf{m}$ . Our objective in this paper is to interpret and assess the reliability of the estimate and its associated uncertainty.

In the Bayesian approach for parameter estimation used by A13,  $\mathbf{m}$  is modeled as a random variable with a specified prior probability distribution that encodes what the authors believed they knew about  $\mathbf{m}$  independent of the measured data. The result of their inference is the posterior probability distribution of  $\mathbf{m}$  conditional on  $\mathbf{d}$ . The likelihood and priors were selected using heuristic arguments that are typical in applications. In particular, they represented their best guesses at interpreting the measured data and a statistical description of their prior knowledge. In this paper we question whether these choices were appropriate, and we describe methods that can be used to validate these choices. We also describe alternative non-Bayesian inversion methods that make different types of assumptions on the statistical model. We will then explain how to compare their results to those from the Bayesian procedures.

The rest of the paper is organized as follows. In section 2 we summarize the work presented in A13. We start the validation analysis by checking the frequentist performance of the Bayesian procedures from A13 in section 3. We focus on the validation of the likelihood and priors in section 4. This analysis leads to a new likelihood and to a revision of the results in A13. The new likelihood is then used in what follows. Since the prior distributions for the experiment are not defined based on the physics of the problem, but on intuitive and convenient choices, it is worth comparing the Bayesian results to those obtained with non-Bayesian methods. In section 5 we find parameter estimates and confidence intervals using maximum likelihood and nonlinear least-squares. For these methods,  $\mathbf{m}$  is no longer modeled as random. In section 6, we use a method that does not require parameter estimates to construct the confidence sets. We summarize and discuss all the results in section 7.



**Fig. 2.1** *Left: PDF (2.1) of the noise  $\varepsilon_i$  plotted as a function of  $\varepsilon_i/\delta_i$ . Right: Function (2.3) that defines the prior density (2.3) for  $t_0$ .*

**2. Summary of A13.** We start by summarizing the Bayesian framework used in A13. There, the noise variables,  $\varepsilon_i = d_i - z_i(\mathbf{m})$ , are assumed to be independent with a shifted/scaled beta(2, 2) distribution: for each  $i$ ,  $\varepsilon_i/\delta_i = 2U_i - 1$ , where the  $U_i$  are independent and identically distributed (i.i.d.) beta(2, 2) random variables. Thus, the probability density function (PDF) of  $\varepsilon_i$  is (see the left panel in Figure 2.1)

$$(2.1) \quad f_i(\varepsilon_i) \propto (1 - \varepsilon_i^2/\delta_i^2) \mathbb{I}_{[-1,1]}(\varepsilon_i/\delta_i),$$

where  $\mathbb{I}_A$  denotes the indicator function of the set  $A$ . For simplicity, normalization constants are omitted here and in what follows. Thus, the PDF of  $\varepsilon_i$  is symmetric about zero with support  $[-\delta_i, \delta_i]$ . Note that  $\delta_i$  is not the standard deviation of  $\varepsilon_i$  but a scaling factor of it:  $\text{Std}(\varepsilon_i) = 4\delta_i^2 \text{Std}(\text{beta}(2, 2))$ . This probability density was chosen by the authors of A13 to indicate that they truly believed that the “real” value should be in the interval  $[d_i - \delta_i, d_i + \delta_i]$ , and within this interval with larger probability in the vicinity of the center than near the boundaries. This PDF leads to the following likelihood of  $\mathbf{m}$ :

$$(2.2) \quad L(\mathbf{m}; \mathbf{d}) \propto \prod_{i=1}^{31} [1 - (d_i - z_i(\mathbf{m}))^2/\delta_i^2] \mathbb{I}_{[-1,1]}([d_i - z_i(\mathbf{m})]/\delta_i).$$

To define the prior distribution for  $\mathbf{m}$ , the parameters  $g$ ,  $C$ , and  $t_0$  are assumed to be independent with  $g$  and  $C$  uniformly distributed on the intervals  $[0, 20]$  (m/s<sup>2</sup>) and  $[0, 0.5]$  (m<sup>-1</sup>), respectively. The authors of A13 chose these as reflecting their best prior knowledge: both parameters certainly need to be positive, and an initial experiment may have provided plausible upper bounds of 20 and 0.5. When  $t_0$  is also modeled as random (the three-parameter case), the authors of A13 chose the prior density for  $t_0$  as shown in the right panel of Figure 2.1 and defined as

$$(2.3) \quad \pi_{t_0}(t_0) \propto \begin{cases} 1, & 41\Delta t \leq t_0 \leq 42\Delta t, \\ 1 - (t_0 - 42\Delta t)/\Delta t, & 42\Delta t \leq t_0 \leq 43\Delta t, \\ 0, & t_0 < 42\Delta t \text{ or } t_0 > 43\Delta t, \end{cases}$$

where  $\Delta t = 1/30$  (s) and  $t$  is measured from the beginning of the recording, predating the time when the box started to fall. Again, this reflected the best heuristic interpretation of the data they obtained from their experiment: (i) they knew that the body

did not move before frame 41; (ii) they could not say for certain whether the body moved between frames 41 and 42; and (iii) it did move for certain in frame 43. In other words,  $\pi_{t_0}$  describes their beliefs about the time  $t_0$  when the body might have started to move.

Given all this, the prior PDF of  $\mathbf{m}$  is either

$$(2.4) \quad \pi(\mathbf{m}) \propto \mathbb{I}_{[0,20]}(g) \mathbb{I}_{[0,0.5]}(C),$$

when  $t_0$  is fixed (the two-parameter case), or

$$(2.5) \quad \pi(\mathbf{m}) \propto \mathbb{I}_{[0,20]}(g) \mathbb{I}_{[0,0.5]}(C) \pi_{t_0}(t_0)$$

in the three-parameter case. In what follows, we will consider these as either separate cases or as two different priors: in the first case, the (very informative) prior distribution for the start time  $t_0$  simply assumes that it is known.

In A13, moments of the posterior density  $f_{\mathbf{m}|\mathbf{d}}(\mathbf{m} | \mathbf{d}) \propto L(\mathbf{m}; \mathbf{d}) \pi(\mathbf{m})$ , in particular the mean and standard deviation, are evaluated by quadrature approximations of the integrals. This was possible because of the low-dimensional nature of the parameter estimation problem.<sup>2</sup> The posterior means of the parameters reported in A13 for priors (2.4) and (2.5) are, respectively,

$$(2.6) \quad \mathbb{E}(g, C | \mathbf{d}) = (8.82, 0.12) \quad \text{and} \quad \mathbb{E}(g, C, t_0 | \mathbf{d}) = (8.64, 0.11, 1.40).$$

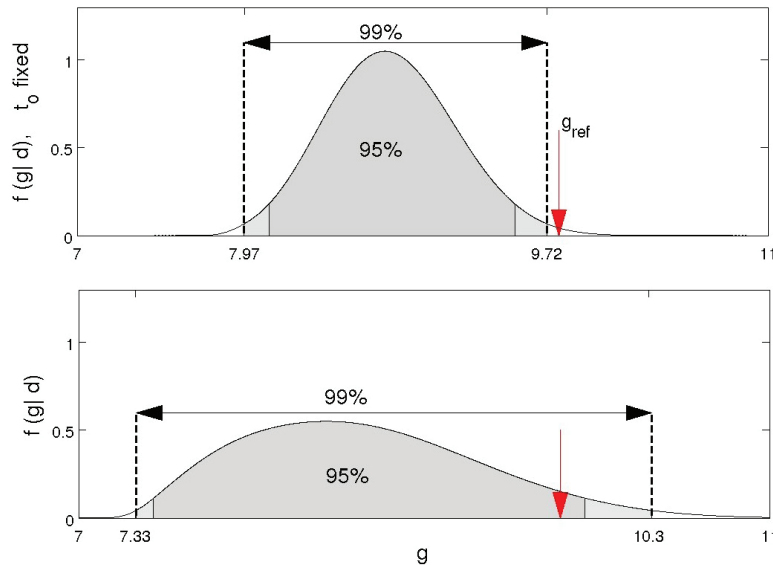
Using the fixed  $t_0$  changes the posterior mean of  $g$  by 2% and that of  $C$  by 9%, but the difference in the results is more obvious in the posterior distributions. The posterior densities for  $g$  using either of the two priors for  $\mathbf{m}$  are shown in Figure 2.2. The figure also shows the highest posterior density (HPD) 95% and 99% credible intervals for  $g$  (in m/s<sup>2</sup>):

	Prior	95%	99%
(2.7a)	(2.4) :	[8.11, 9.53]	[7.97, 9.72]
(2.7b)	(2.5) :	[7.43, 9.93]	[7.33, 10.31].

Recall that credible intervals are not uniquely defined and, in particular, they can have very different lengths even if they have the same level  $1 - \alpha$  (as usual, we refer to the level of credible or confidence sets as  $1 - \alpha$  for some  $\alpha \in (0, 1)$ ). When possible, HPD intervals are used as they provide intervals of minimum length for a fixed level  $1 - \alpha$  [3]. We note that without the prior on  $t_0$  (i.e., for the two-parameter case), neither the 99% nor the 95% credible interval contains  $g_{\text{ref}}$  and  $\mathbb{P}(g \geq g_{\text{ref}} | \mathbf{d}) \approx 0.006$ . Hence,  $g_{\text{ref}}$  seems to be an unlikely value under this posterior distribution. In contrast, the corresponding intervals both contain  $g_{\text{ref}}$  when using the prior on  $t_0$ . This prior adds variability to the posterior, which makes them right-skewed and yields intervals that are wide enough to include  $g_{\text{ref}}$ .

Throughout the rest of the paper, Bayesian results will be compared primarily through their posterior means,  $\mathbb{E}(g, C | \mathbf{d})$  or  $\mathbb{E}(g, C, t_0 | \mathbf{d})$ , and credible intervals. To compare these results to those obtained with non-Bayesian methods, we will use the frequentist coverage of credible intervals, which will be defined in section 3. Results from the various methods we consider are summarized and discussed in section 7.

<sup>2</sup>In more complex problems, one must approximate integrals using, for example, Markov chain Monte Carlo sampling, and in that case one would also need to check convergence of Markov chains in the validation process.



**Fig. 2.2** *Top: Posterior of  $g$  corresponding to likelihood (2.2) and prior (2.4). The shaded areas define the 95% and 99% HPD credible intervals for  $g$ . Bottom: Corresponding results for prior (2.5). The value  $g_{\text{ref}}$  (1.2) is shown for reference and marked by the vertical arrows.*

Let us now conclude our summary of A13. There, the final result was the posterior density for  $g$  under prior (2.5), and corresponding credible intervals. One of the points of this paper is a validation analysis. Such checks are necessary to make sure the statistical model used is appropriate and the estimates are meaningful.

**3. A First Validation Check.** Assume for the moment that the likelihood function (2.2) is appropriate. If we asked the Texas A&M team to repeat the experiment 100 times with exactly the same setup, and if for each experiment they constructed 95% credible intervals for  $g$  as above, would they find that approximately 95 of the intervals contain the “true”  $g_{\text{ref}}$ ? In statistics this is known as a check of the frequentist coverage of the credible intervals constructed above. If the frequentist coverage matches the target 95%, then we may conclude that despite the ad-hoc choice of priors, the Bayesian credible intervals can be interpreted as frequentist 95% confidence intervals. This provides one possible interpretation of the posterior credible intervals, and is reassuring to those who are concerned with subjectivity in the selection of the priors. For further discussion and examples of frequentist performance of Bayesian procedures, see [3, 10, 20].

We use simulations to obtain information about the frequentist coverage and mean length of the credible intervals using the procedure described in Algorithm 1. The basic idea is to generate synthetic data from a fixed  $\widehat{\mathbf{m}}$ , which we assume to be the “true”  $\mathbf{m}$ , and to repeat the estimation process many times. This method is known as parametric bootstrap [5, 8] because the noise is sampled from a parametric family (a beta distribution in this case). If the likelihood is correct, the simulations can be thought of as synthetic repetitions of the experiment with a fixed known  $\mathbf{m}$ .

The results of the simulations for 95% credible intervals are shown in Table 3.1. The table shows the usual 95% Agresti–Coull intervals [2] for the coverage and 95%

**Algorithm 1** Simulations to estimate the frequentist coverage and mean length of  $1 - \alpha$  intervals  $I_\alpha(g)$  for  $g$ .

**fix:**  $\mathbf{m} = \widehat{\mathbf{m}}$ , or  $\mathbf{m} = (g_{\text{ref}}, \widehat{C}, t_0)$ , or  $\mathbf{m} = (g_{\text{ref}}, \widehat{C}, \widehat{t}_0)$  as the “true” parameters  
**for**  $k = 1 : N$  **do**  
 generate:  $U_1, \dots, U_{31}$  i.i.d. beta(2, 2)  
 define  $\varepsilon_i^* = \delta_i(2U_i - 1)$ ,  $\mathbf{d}^* = \mathbf{z}(\mathbf{m}) + \varepsilon^*$   
 construct  $1 - \alpha$  credible interval  $I_\alpha(g)$  for  $g$  by evaluating  $f_{\mathbf{m}|\mathbf{d}}(\mathbf{m}|\mathbf{d}^*)$   
**end for**  
 Frequentist coverage is estimated by the proportion of intervals that contain the “true”  $g$ .

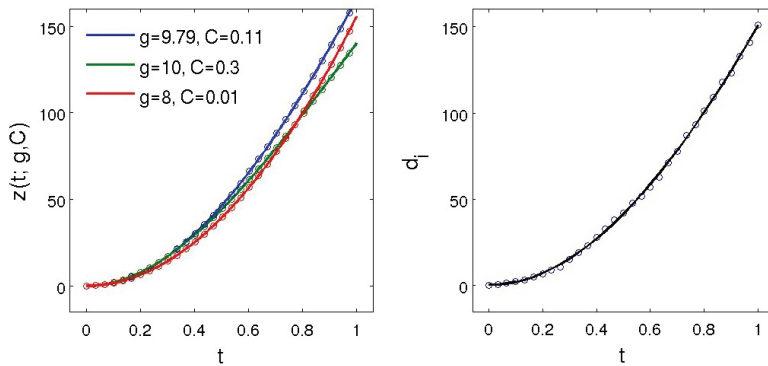
**Table 3.1** Frequentist coverage and mean length of the 95% credible intervals based on likelihood (2.2) and the two different priors on  $\mathbf{m}$ , (2.4) and (2.5).

$\pi(\mathbf{m})$	$\widehat{g} = \mathbb{E}(g   \mathbf{d})$		$\widehat{g} = g_{\text{ref}}$	
	Coverage	Average length	Coverage	Average length
(2.4)	95.13 ± 0.01	0.83 ± 0.01	94.63 ± 0.01	0.84 ± 0.01
(2.5)	97.12 ± 0.01	2.03 ± 0.02	97.02 ± 0.01	2.20 ± 0.02

large-sample confidence intervals for the mean interval length [16]. We see that the intervals based on the prior with  $t_0$  fixed have approximately 95% frequentist coverage but those based on the prior (2.5) are much wider and more conservative. Judging from this frequentist performance, we may question whether the prior for  $t_0$  adds variability (thus making the intervals wide enough to cover  $g_{\text{ref}}$ ) that may not be justified. We will address this question in the next section by checking the likelihood function.

**4. Model Validation.** The results of the previous section have revealed possible problems with the credible intervals for three parameters as derived in A13. We will show in this section that the problem stems from an overestimation of the measurement uncertainties.

**4.1. Checking the Noise Level.** We start by validating the assumed distribution of the noise whose PDF is given by (2.1). As mentioned before, this probability density was chosen by physical intuition, but is the distribution of the measurement noise assumed by the authors of A13 appropriate? The idea of this section is to “get to the noise”  $\varepsilon$  by somehow subtracting  $\mathbf{z}(\mathbf{m})$  from  $\mathbf{d}$ . Since in this case estimating  $\mathbf{m}$  is not the goal, we may take advantage of a surrogate model for  $\mathbf{z}(\mathbf{m})$  that leads to residuals that are easier to analyze. We first note that a cubic fit to  $z_i(\mathbf{m})$  (i.e., polynomial regression [6]) provides a good approximation. For example, the left panel in Figure 4.1 shows values of  $\log \cosh [\sqrt{gC} t] / C$  (circles) at different times, and its cubic fit (lines) for three different choices of  $g$  and  $C$ . In each case there is a reasonably good agreement (the maximum difference is of order  $10^{-3}$ ). Therefore, to study the noise distribution, we model the data as  $\mathbf{d}_\delta = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_\delta$ , where  $\mathbf{d}_\delta = \mathbf{D}\mathbf{d}$ ,  $\boldsymbol{\varepsilon}_\delta = \mathbf{D}\boldsymbol{\varepsilon}$ ,  $\mathbf{D} = \text{diag}\{1/\delta_1, \dots, 1/\delta_{31}\}$ , and the  $31 \times 4$  matrix  $\mathbf{X}$  has columns  $(1/\delta_i)$ ,  $(t_i/\delta_i)$ ,  $(t_i^2/\delta_i)$ , and  $(t_i^3/\delta_i)$ . Note that  $\boldsymbol{\varepsilon}_\delta$  is unit-less. The elements of the vector  $\boldsymbol{\beta}$  we want to estimate are then the expansion coefficients of the cubic polynomial we want to fit to the data. The entries of the vector of differences between data and cubic model,



**Fig. 4.1** Left: Function  $z(\mathbf{m})$  (1.3) as a function of time for three different choices of  $g$  and  $C$ . The lines show the cubic approximations. Right: Cubic fit (line) to the data (circles).

$\varepsilon_\delta$ , are independent, zero-mean, with

$$(4.1) \quad (\varepsilon_\delta)_i = 2U_i - 1, \quad U_i \text{ i.i.d. beta}(2, 2).$$

We use least-squares to estimate  $\beta$ : we find the vector  $\hat{\beta}$  that minimizes  $\|\mathbf{d}_\delta - \mathbf{X}\beta\|^2$  over  $\beta \in \mathbb{R}^4$ , in other words, a weighted least-squares fit to  $\mathbf{d}$ , where each point is weighted by  $1/\delta_i$ . The fit is shown in the right panel of Figure 4.1. The residual vector of the fit,  $\mathbf{r} = \mathbf{d}_\delta - \mathbf{X}\hat{\beta}$ , has zero mean and covariance matrix  $\Sigma_{\mathbf{r}} = \sigma^2 \mathbf{Q}$ , where  $\sigma^2 = \text{Var}(2U) = 1/5$  and  $\mathbf{Q}$  is the projection matrix onto the orthogonal complement of the column space of  $\mathbf{X}$  (see, for example, [6]):  $\mathbf{Q} = \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .

We transform the residuals to correct for their correlation and heterogeneous variances. Since  $\mathbf{Q}$  is a projection matrix with a four-dimensional nullspace, we can write its SVD as

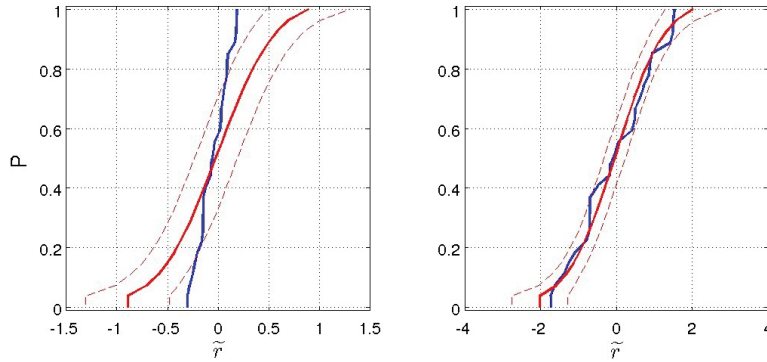
$$\mathbf{Q} = \mathbf{U} \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}^\top = (\mathbf{U}_1 \quad \mathbf{U}_2) \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{U}_1^\top \\ \mathbf{U}_2^\top \end{pmatrix} = \mathbf{U}_1 \mathbf{U}_1^\top,$$

where  $\mathbf{U}$  is an orthogonal matrix,  $\mathbf{I}$  is the  $27 \times 27$  identity matrix, and  $\mathbf{U}_1$  and  $\mathbf{U}_2$  have orthonormal columns. Furthermore, since  $\mathbf{r} = \mathbf{Q}\varepsilon_\delta$ , it follows that

$$(4.2) \quad \tilde{\mathbf{r}} \equiv \mathbf{U}_1^\top \mathbf{r} = \tilde{\varepsilon} \equiv \mathbf{U}_1^\top \varepsilon_\delta.$$

In particular,  $\tilde{\mathbf{r}}$  and  $\tilde{\varepsilon}$  have the same distribution, which helps us check the posited distribution of  $\varepsilon$ . To compare distributions, we will use empirical cumulative distribution functions (CDFs), which are often used for validation and do not require the choice of bins as do histograms (see, for example, [16]). We simulate  $\text{nsims}=10,000$  noise vectors  $\varepsilon_\delta$  with distribution (4.1). For each realization of  $\varepsilon_\delta$ , we compute the empirical CDF,  $\hat{F}_{\tilde{\varepsilon}}(x)$ , of  $\tilde{\varepsilon}$  on a uniform grid of values,  $x$ , of the random variable. The solid red line in the left panel in Figure 4.2 is the sample mean approximation of  $\mathbb{E}[\hat{F}_{\tilde{\varepsilon}}(x)]$  (i.e.,  $\text{mean}(\hat{F}_{\tilde{\varepsilon}}(x))$ ), and the dashed lines define 95% confidence intervals for  $\mathbb{E}[\hat{F}_{\tilde{\varepsilon}}(x)]$  given by  $\text{mean}(\hat{F}_{\tilde{\varepsilon}}(x)) \pm 1.96 \text{std}(\hat{F}_{\tilde{\varepsilon}}(x))/\sqrt{\text{nsims}}$ , where  $\text{mean}$  and  $\text{std}$  denote the usual sample mean and sample standard deviation (over the  $\text{nsims}$  simulations). The empirical CDF of  $\tilde{\mathbf{r}}$  is shown in blue. The plot shows that the assumed distribution of the noise is not consistent with that of the noise in the data: the true measurement noise appears to be much more narrowly centered around zero than the





**Fig. 4.2** *Left: Empirical CDF of  $\tilde{r}$  (blue) and the mean and standard deviation of the empirical CDF of  $\tilde{\epsilon}$  (red) when the noise is assumed to have PDF (2.1). Right: Empirical CDF of  $\tilde{r}$  (blue) and empirical CDF of a standardized sequence of Gaussian random variables (red).*

modeled uncertainty described by the noise PDF. In fact, even the standard deviations are not consistent: our estimate of the standard deviation of  $d_i$  is  $0.15\delta_i$  as compared to the assumed  $\sqrt{1/5}\delta_i \approx 0.45\delta_i$ . It turns out that we can use a Gaussian distribution to model these data. To see this, note that if  $\epsilon$  is multivariate Gaussian, then so are  $\epsilon_\delta$  and  $U_1^\top \epsilon_\delta$ . Moreover, the entries of  $U_1^\top \epsilon_\delta$  are i.i.d. Gaussian because the columns of  $U_1$  are orthonormal. This means that the entries of  $\tilde{r}$  are also i.i.d. Gaussian and therefore each of the corrected entries of  $\tilde{r}$ ,  $T_i = (\tilde{r}_i - \text{mean}(\tilde{r}))/\text{std}(\tilde{r})$ , has a Student  $t_{26}$  distribution. Since these centered residuals are correlated, we conduct simulations to estimate the pointwise mean and standard deviation of the empirical CDF of 27 i.i.d. Gaussian variables that are corrected by subtracting their sample mean and dividing by their sample standard deviation. The right panel in Figure 4.2 shows the results. The blue line corresponds to the empirical CDF of the variables  $T_i$ . We see that this Gaussian distribution seems more reasonable than the scaled/shifted beta(2, 2) distribution. Other distributions could be used but the Gaussian is convenient, especially for the estimation methods we will use below.

The primary conclusion of this analysis is that while the authors of A13 might have had good reasons to choose their noise model (2.1) with measured values for  $\delta_i$ , the noise in the data does not match this assumption; the true distribution has a smaller standard deviation. In other words, their distributional choice was overly cautious, a fact already acknowledged in A13. The consequence is that their credible intervals are unnecessarily large and not backed up by the actual data.

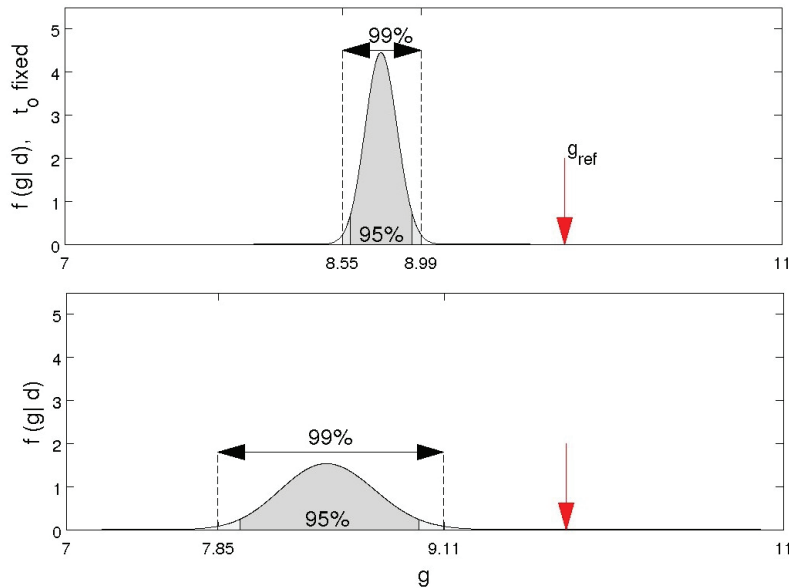
The analysis above leads us to assume instead that the variables  $\epsilon_i/\delta_i$  are i.i.d.  $N(0, \sigma_{\text{cubic}}^2)$ , where  $\sigma_{\text{cubic}}^2$  is the standard least-squares estimate of  $\sigma^2$  (unit-less) obtained from the cubic fit. Thus, the likelihood is now

$$(4.3) \quad L(\mathbf{m}; \mathbf{d}) \propto \prod_{i=1}^{31} \exp \left[ -\frac{(d_i - z_i(\mathbf{m}))^2}{2\delta_i^2 \sigma_{\text{cubic}}^2} \right],$$

and we can repeat the analysis to obtain

$$(4.4) \quad \mathbb{E}(g, C | \mathbf{d}) = (8.77, 0.11) \quad \text{and} \quad \mathbb{E}(g, C, t_0 | \mathbf{d}) = (8.47, 0.09, 1.39).$$

The largest difference between these results and those of A13 (stated in (2.6)) is in the posterior mean of  $g$ . This is more evident when we compare the posterior distributions.



**Fig. 4.3** Same as Figure 2.2 but using the Gaussian likelihood (4.3).

**Table 4.1** Same as Table 3.1 but using the Gaussian likelihood (4.3).

$\pi(\mathbf{m})$	$\hat{g} = \mathbb{E}(g   \mathbf{d})$		$\hat{g} = g_{\text{ref}}$	
	Coverage	Average length	Coverage	Average length
(2.4)	$95.03 \pm 0.01$	$0.34 \pm 0.01$	$95.13 \pm 0.01$	$0.35 \pm 0.01$
(2.5)	$94.13 \pm 0.01$	$0.99 \pm 0.01$	$94.23 \pm 0.01$	$1.01 \pm 0.01$

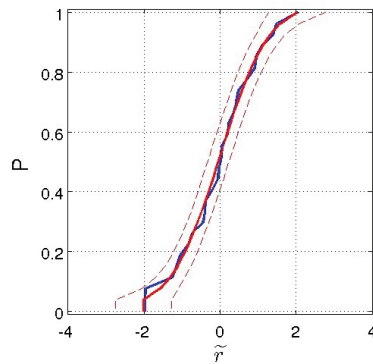
Figure 4.3 shows that the new posteriors are indeed tighter and more symmetric than before, and that  $g_{\text{ref}}$  is an even more unlikely value than under likelihood (2.2), whether or not we use the prior on  $t_0$ . The 95% and 99% HPD credible intervals for  $g$  are now (compare these with (2.7))

$$(4.5a) \quad \begin{array}{l} \text{Prior} \\ (2.4) : \end{array} \quad \begin{array}{cc} 95\% & 99\% \\ [8.60, 8.94] & [8.55, 8.99] \end{array}$$

$$(4.5b) \quad \begin{array}{l} \text{Prior} \\ (2.5) : \end{array} \quad \begin{array}{cc} 95\% & 99\% \\ [7.97, 8.97] & [7.85, 9.11]. \end{array}$$

The frequentist coverage of these intervals is given in Table 4.1. The new coverage for prior (2.4) is closer to 95% and it is achieved with intervals that are, on average, less than half the size. The coverage for (2.5) is also closer to the target 95%, again with intervals half the size. This shows that the new likelihood, in addition to being more consistent with the noise in the data, leads to Bayesian results that have good frequentist performance.

**4.2. Reconsidering the Use of Individual Uncertainty Measures  $\delta_i$ .** It may be intuitively appealing to control the uncertainty at every point with the uncertainty factors  $\delta_i$ , as is done in A13, but is it necessary? One can argue that these uncertainties reflect the rotation of the falling box and errors caused by the resolution of the camera,



**Fig. 4.4** Same as in the right panel in Figure 4.2 but using the likelihood (4.6) without  $\delta_i$ .

but we believe that the uncertainty in  $d_i$  may not change much in time and that it may not be necessary to include different  $\delta_i$  in the likelihood. In this case, we could compute the cubic approximation  $\beta$  weighing all points equally. Figure 4.4 is the equivalent of the right panel in Figure 4.2 but using this equally weighted cubic fit. The results seem to indicate that the residuals are also consistent with the noise when we use the likelihood with equal  $\delta_i$ . In fact, this plot looks better than the one that uses the  $\delta_i$ . In the spirit of using the simplest model consistent with the data, we shall assume henceforth that the errors  $\varepsilon_i$  are i.i.d.  $N(0, \sigma_{\text{cubic}}^2)$ , where  $\sigma_{\text{cubic}}^2$  is now obtained from the cubic fit without  $\delta_i$  (this time  $\sigma_{\text{cubic}}$  has the same units as  $d_i$ ). Thus, the likelihood that will be used in the rest of the paper is:

$$(4.6) \quad L(\mathbf{m}; \mathbf{d}) \propto \prod_{i=1}^{31} \exp \left[ -\frac{(d_i - z_i(\mathbf{m}))^2}{2\sigma_{\text{cubic}}^2} \right],$$

which leads to the posterior means

$$(4.7) \quad \mathbb{E}(g, C) | \mathbf{d} = (8.77, 0.11) \quad \text{and} \quad \mathbb{E}(g, C, t_0) | \mathbf{d} = (8.47, 0.09, 1.39).$$

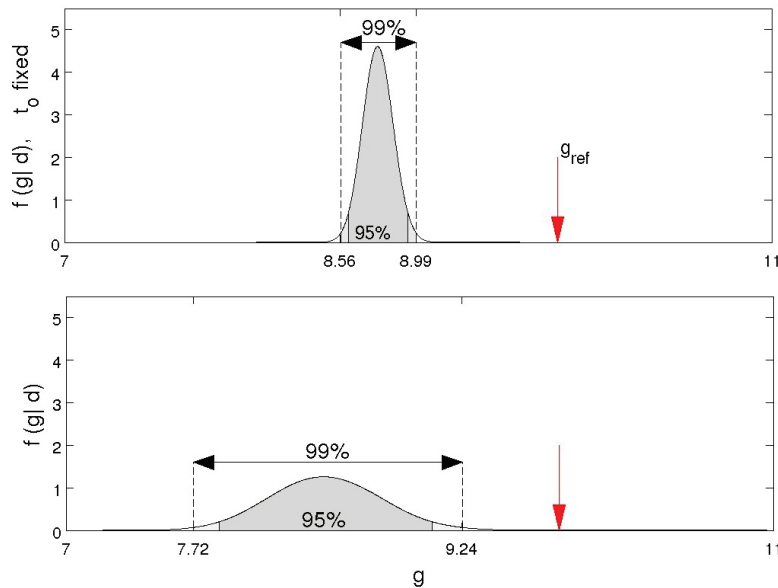
These are almost the same as (4.4), which were obtained using  $\delta_i$ . The corresponding posterior distributions are shown in Figure 4.5. The new credible intervals for  $g$  are

$$(4.8a) \quad \begin{array}{ccc} \text{Prior} & 95\% & 99\% \\ (2.4) & : [8.61, 8.94] & [8.56, 8.99] \end{array}$$

$$(4.8b) \quad (2.5) : [7.87, 9.07] \quad [7.72, 9.24].$$

There are some small differences with (4.5), but to compare them, it is better to look at their frequentist coverage and mean length. The estimated frequentist coverage of the 95% credible intervals is shown in Table 4.2. The coverage and mean length of the intervals under prior (2.4) are similar to those that used different  $\delta_i$ . The coverage under (2.5) is a little above the target (compared to slightly below with the  $\delta_i$ ) and the intervals are slightly wider. However, overall, the use of  $\delta_i$  seems to make little difference.

**4.3. Prior and Posterior Predictive Checks.** If the likelihood function and prior distribution provide an appropriate model for the data, we should be able to generate

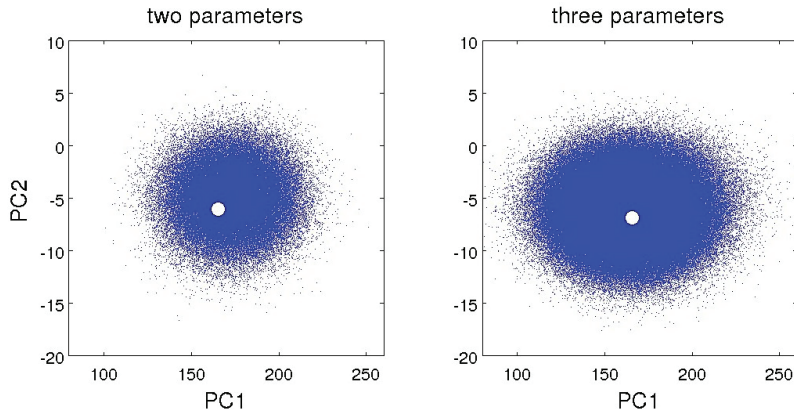


**Fig. 4.5** Same as Figure 2.2 but using the likelihood (4.6) that does not use  $\delta_i$ .

**Table 4.2** Same as Table 3.1 but using the likelihood (4.6) that does not use  $\delta_i$ .

$\pi(\mathbf{m})$	$\hat{g} = \mathbb{E}(g   \mathbf{d})$		$\hat{g} = g_{\text{ref}}$	
	Coverage	Average length	Coverage	Average length
(2.4)	$95.23 \pm 0.01$	$0.33 \pm 0.01$	$95.03 \pm 0.01$	$0.34 \pm 0.01$
(2.5)	$95.72 \pm 0.01$	$1.18 \pm 0.01$	$95.53 \pm 0.01$	$1.23 \pm 0.01$

data that are similar to the original  $\mathbf{d}$ . For example, we can simulate data  $\mathcal{D}_{\text{pri}} = \{\mathbf{d}_1^*, \dots, \mathbf{d}_m^*\}$  by drawing  $m$  samples  $\mathbf{m}_1^*, \dots, \mathbf{m}_m^*$  from the prior distribution and then generating  $\mathbf{d}_i^*$  from the conditional distribution  $F_{\mathbf{d}|\mathbf{m}}(\mathbf{d} | \mathbf{m}_i^*)$  using the forward model (1.1). The distribution of the samples in  $\mathcal{D}_{\text{pri}}$  is called a prior predictive distribution. If instead we sample  $\mathbf{m}_i^*$  from the posterior distribution of  $\mathbf{m}$  conditional on  $\mathbf{d}_i^*$ , then the distribution of the sample,  $\mathcal{D}_{\text{post}}$ , is called a posterior predictive distribution [3, 13]. If the prior and likelihood are reasonable, we would expect the characteristics of  $\mathbf{d}$  to be similar to those of the samples in  $\mathcal{D}_{\text{pri}}$  or  $\mathcal{D}_{\text{post}}$ . In particular, if  $\hat{\boldsymbol{\beta}}$  are the coefficients of the cubic fit to  $\mathbf{d}$  and  $\hat{\boldsymbol{\beta}}_i^*$  those of the cubic fit to  $\mathbf{d}_i^*$ , then  $\boldsymbol{\beta}$  should be consistent with the distribution of the  $\hat{\boldsymbol{\beta}}_i^*$ . To check this we can, for example, compute the principal components (PCs) [14] of the data matrix  $[\hat{\boldsymbol{\beta}}_1^* \dots \hat{\boldsymbol{\beta}}_m^*]$  and plot the first two coefficients of  $\hat{\boldsymbol{\beta}}$  (white circle) and  $\hat{\boldsymbol{\beta}}_i^*$  (blue) with respect to the PCs (which account for more than 98% of the total variability). To sample from the posterior predictive distributions, we have used the Metropolis–Hastings algorithm [13]. Figure 4.6 shows the results. We would question the choice of priors if the white circles were far from the bulk of the blue points. In this case, we see that the PC coefficients of  $\hat{\boldsymbol{\beta}}$  are consistent with realizations of the PC coefficients of  $\hat{\boldsymbol{\beta}}^*$  under the assumed model with two or three parameters. The same happens for the prior predictive distributions.



**Fig. 4.6** Left: Plot of  $\beta$  (white) and  $\beta_i^*$  (blue) in the first two PCs for prior (2.4). Right: The same for prior (2.5).

To provide an example, we have used only one function of the data to check the priors. Usually we try to use functions that reduce the dimensionality of  $\mathbf{d}$ ; it was reduced to two using the first two PCs. Many other functions,  $T(\mathbf{d})$ , can be used to compare to characteristics of  $T(\mathbf{d}_i^*)$ . Of course, we could find functions for which we see some discrepancies, but the idea is to choose a function that is relevant to the problem. More examples on the use of predictive distributions for validation can be found in [3, 13], and also in [10, 21] in the context of inverse problems.

**5. Maximum Likelihood and Nonlinear Least-Squares.** It is important to know how to use a variety of estimation methods that make different assumptions and incur different computational costs. Since different methods used soundly should lead to results that are consistent, comparing their results is one way to look for possible problems. Up to now we have only considered the Bayesian inference of  $\mathbf{m}$ . This parameter vector was modeled as random with a prior probability distribution, and the results of the inference were conditional on the data. But parameter estimates can also be obtained without assuming that  $\mathbf{m}$  is random. Inferential methods in which the unknown parameter is fixed (not random) and whose performance is based on their distribution induced by different realizations of the data are called frequentist. An introduction to Bayesian and non-Bayesian methods for inverse problems can be found in [10]. A more general reference is [17].

**5.1. Maximum Likelihood Estimates.** Let us first consider the method of maximum likelihood (ML) to obtain parameter estimates without using prior distributions. It is based only on  $\mathbf{d}$  and the corresponding likelihood function defined by (4.6). A value,  $\widehat{\mathbf{m}}_{\text{ML}}$ , that maximizes the likelihood,  $L(\mathbf{m}; \mathbf{d})$ , is called an ML estimate of  $\mathbf{m}$ . To find  $\widehat{\mathbf{m}}_{\text{ML}}$ , we evaluate  $L(\mathbf{m}; \mathbf{d})$  on a grid of parameter values.<sup>3</sup> The ML estimates for the two cases  $\mathbf{m} = (g, C)$  and  $\mathbf{m} = (g, C, t_0)$  are, respectively,

$$(5.1) \quad (\widehat{g}, \widehat{C})_{\text{ML}} = (8.77, 0.11) \quad \text{and} \quad (\widehat{g}, \widehat{C}, \widehat{t}_0)_{\text{ML}} = (8.41, 0.09, 1.39).$$

We observe that these ML estimates are almost identical to the corresponding posterior means under priors (2.4) and (2.5). This is not surprising given the flat priors for  $g$  and  $C$  and the “almost” flat prior for  $t_0$ .

<sup>3</sup>For consistency, we use the same grid of parameter values used to compute the integrals in the Bayesian calculations conducted by A13.

In addition to parameter estimates, ML can be used to obtain approximate confidence intervals, which are derived using the asymptotic theory of ML estimators. Under regularity conditions, ML estimators are asymptotically unbiased and Gaussian with covariance matrix equal to the inverse Fisher information matrix [4, 9]. The Fisher information matrix of a  $k \times 1$  parameter vector  $\mathbf{m}$  is the  $k \times k$  matrix with entries

$$(5.2) \quad \mathbf{I}(\mathbf{m})_{ij} = \mathbb{E} \left[ \frac{\partial^2 (-\ln L(\mathbf{m}; \mathbf{d}))}{\partial m_i \partial m_j} \right].$$

An approximate  $1 - \alpha$  confidence interval for  $m_j$  is given by [4]

$$(5.3) \quad (\widehat{\mathbf{m}}_{\text{ML}})_j \pm z_{\alpha/2} \sqrt{[\mathbf{I}(\widehat{\mathbf{m}}_{\text{ML}})^{-1}]_{jj}},$$

where  $z_{\alpha/2}$  is the usual  $1 - \alpha/2$  quantile of the  $N(0, 1)$  distribution. When the expectation in (5.2) cannot be easily computed, it may be approximated using sample averages. This leads to the observed Fisher information matrix  $\mathbf{J}(\mathbf{m})$ ,

$$(5.4) \quad \mathbf{J}(\mathbf{m})_{ij} = \sum_{\ell=1}^K \frac{\partial^2 (-\ln L(\mathbf{m}; d_\ell))}{\partial m_i \partial m_j},$$

and approximate  $1 - \alpha$  confidence intervals  $(\widehat{\mathbf{m}}_{\text{ML}})_j \pm z_{\alpha/2} \sqrt{[\mathbf{J}(\widehat{\mathbf{m}}_{\text{ML}})^{-1}]_{jj}}$  [7]. Since we are using a Gaussian likelihood, it is easy enough to use  $\mathbf{I}(\widehat{\mathbf{m}}_{\text{ML}})$  instead of  $\mathbf{J}(\widehat{\mathbf{m}}_{\text{ML}})$ . The approximate 95% and 99% confidence intervals (5.3) for  $g$  for the two- and three-parameter cases are

$$(5.5a) \quad \mathbf{m} \quad 95\% \quad 99\%$$

$$(g, C) : [8.60, 8.94] \quad [8.55, 8.99]$$

$$(5.5b) \quad (g, C, t_0) : [7.78, 9.03] \quad [7.59, 9.23].$$

The intervals for the two-parameter case are almost identical to the corresponding credible intervals (4.8). The intervals for the three-parameter case seem to be slightly wider than the corresponding credible intervals, but since the two types of intervals have different interpretations, we compare them through their frequentist coverage and mean lengths. Also, since the ML intervals are based on asymptotic approximations, we should assess their actual coverage with  $n = 31$ . To this end, we conduct simulations using the following slight modification of Algorithm 1: In the construction step, we compute ML intervals using (5.3) instead of credible intervals. The results are shown in Table 5.1. We see that the coverage and mean length of the ML intervals and credible intervals are very similar, although the latter seem to be slightly shorter.

**Table 5.1** Interval length and coverage results for the approximate 95% ML confidence intervals using the Gaussian likelihood (4.6) with two or three parameters.

$\mathbf{m}$	$\widehat{g} = \widehat{g}_{\text{ML}}$		$\widehat{g} = g_{\text{ref}}$	
	Coverage	Average length	Coverage	Average length
$(g, C)$	$94.83 \pm 0.01$	$0.34 \pm 0.01$	$94.33 \pm 0.01$	$0.35 \pm 0.01$
$(g, C, t_0)$	$94.93 \pm 0.02$	$1.26 \pm 0.01$	$95.15 \pm 0.01$	$1.29 \pm 0.01$

**5.2. Nonlinear Least-Squares.** One advantage of least-squares (LS), linear or nonlinear, is that no distributional assumptions are required to obtain the estimates. We need only find a vector of parameters,  $\widehat{\mathbf{m}}_{\text{NR}}$ , that minimizes  $\|\mathbf{d} - \mathbf{z}(\mathbf{m})\|^2$ . Knowledge of the noise distribution is needed, however, to construct confidence intervals. In the nonlinear Gaussian case, confidence sets are defined by applying the Gaussian theory to a linearization around the estimate. In the case when  $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ , the distribution of  $\widehat{\mathbf{m}}_{\text{NR}}$  is approximately  $N(\mathbf{m}, \sigma^2 [\mathbf{F}(\mathbf{m})^\top \mathbf{F}(\mathbf{m})]^{-1})$ , where  $\mathbf{F}(\mathbf{m})$  is the Jacobian matrix of the forward model with entries  $\mathbf{F}(\mathbf{m})_{ij} = \partial z_i(\mathbf{m}) / \partial m_j$ ; see [18]. An approximate  $1 - \alpha$  confidence interval for  $m_j$  is given by

$$(5.6) \quad (\widehat{\mathbf{m}}_{\text{NR}})_j \pm t_{31-p}(\alpha/2) \hat{\sigma} \sqrt{[(\mathbf{F}(\widehat{\mathbf{m}}_{\text{NR}})^\top \mathbf{F}(\widehat{\mathbf{m}}_{\text{NR}}))^{-1}]_{jj}},$$

where the estimate of the variance,  $\sigma^2$ , is  $\hat{\sigma}^2 = \|\mathbf{d} - \mathbf{z}(\widehat{\mathbf{m}}_{\text{NR}})\|^2 / (31 - p)$  with  $p = 2$  or  $p = 3$ , depending on whether the fit is done to two or three parameters. We again use Algorithm 1 to assess the actual coverage of these approximate confidence intervals (using (5.6) in the construction step).

In the Gaussian likelihood case, ML and nonlinear LS estimates should be exactly the same. In our case there are small differences because the ML estimates were obtained using a grid (so as to be able to compare to the posterior results in A13), while the nonlinear LS fit was done using Levenberg–Marquadt search method [18]. The square root in (5.6) should also be the same as that in (5.3), but for the nonlinear LS estimates used here,  $\sigma^2$  is estimated each time using the residuals of the nonlinear fit, not the residuals of the cubic fit to the surrogate model. Therefore, we expect small differences in the coverage and length of the confidence intervals. The reason for using these different procedures is to check consistency of the results under different but reasonable estimation methods.

The nonlinear LS estimates for the two- and three-parameter cases are

$$(5.7) \quad (\widehat{g}, \widehat{C}) = (8.77, 0.11) \quad \text{and} \quad (\widehat{g}, \widehat{C}, \widehat{t}_0) = (8.45, 0.09, 1.39).$$

The 95% and 99% confidence intervals for  $g$  are given by

$$(5.8a) \quad \begin{array}{ccc} \mathbf{m} & 95\% & 99\% \\ (g, C) & : [8.60, 8.94] & [8.54, 9.01] \end{array}$$

$$(5.8b) \quad \begin{array}{ccc} \mathbf{m} & 95\% & 99\% \\ (g, C, t_0) & : [7.80, 9.10] & [7.57, 9.33]. \end{array}$$

Table 5.2 shows their estimate coverages and mean lengths. Again we get results consistent with ML and with the credible intervals. However, the nonlinear LS intervals tend to be wider than those from the other two methods. This is not surprising, as  $\sigma$  is fixed in the latter while it is estimated in each nonlinear fit.

If we had no idea about the distribution of the noise, we could still use nonlinear LS and we could find approximate confidence intervals using nonparametric bootstrap

**Table 5.2** Interval length and coverage results for the 95% nonlinear LS confidence intervals of  $g$  using likelihood (4.6) with two or three parameters.

$\mathbf{m}$	$\widehat{g} = \widehat{g}_{\text{NR}}$		$\widehat{g} = g_{\text{ref}}$	
	Coverage	Average length	Coverage	Average length
$(g, C)$	$94.33 \pm 0.01$	$0.35 \pm 0.01$	$94.33 \pm 0.01$	$0.36 \pm 0.01$
$(g, C, t_0)$	$95.23 \pm 0.01$	$1.31 \pm 0.01$	$95.13 \pm 0.01$	$1.34 \pm 0.01$

methods. For example, in Algorithm 1, we can draw noise samples by sampling with replacements from a corrected version of the residuals of the surrogate model [5, 8, 10].

We can summarize this section by observing that both ML and nonlinear LS estimates and confidence intervals are very similar to those obtained in previous sections. We take this as additional validation that our results are consistent.

**6. Confidence Intervals without Inversion.** The confidence intervals used in section 5 are of the form

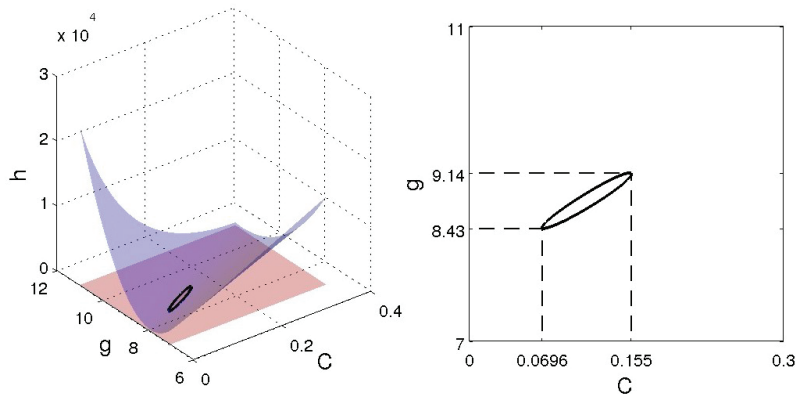
$$\text{estimate} \pm k_\alpha \times (\text{approximate standard error of estimate}),$$

where the constant  $k_\alpha$  is chosen to approximately achieve the desired confidence level  $1 - \alpha$ . This requires solving the inverse problem. In this section we describe a method to construct approximate confidence intervals directly, without first finding an estimate,  $\widehat{\mathbf{m}}$ , and its standard error. The basic tool is a duality between hypothesis tests and confidence sets, which allows us to find a confidence set by inverting a test (e.g., [4, 15]). Some applications of this method to inverse problems can be found in [19].

For the likelihood (4.6), and under the null hypothesis that  $\mathbf{m}$  is the true vector of parameters that generated  $\mathbf{d}$ , we have a  $\chi^2$  random variable with 31 degrees of freedom:  $\|\mathbf{d} - \mathbf{z}(\mathbf{m})\|^2 / \sigma^2 \sim \chi_{31}^2$ . Therefore, the set

$$A_\alpha(\mathbf{m}) = \{\mathbf{d} : \|\mathbf{d} - \mathbf{z}(\mathbf{m})\|^2 \leq \sigma^2 \chi_{31}^2(\alpha)\}$$

defines an acceptance region of significance level  $\alpha$ . A  $1 - \alpha$  confidence region,  $C_\alpha(\mathbf{d})$ , for  $\mathbf{m}$  is defined by inverting the test:  $C_\alpha(\mathbf{d}) = \{\mathbf{m} : \mathbf{d} \in A_\alpha(\mathbf{m})\}$ . We use  $\sigma_{\text{cubic}}$  from section 4.2 to approximate  $\sigma$ . Figure 6.1 shows the boundary of the confidence region for the two-parameter case and  $\alpha = 0.05$ . It is approximately an ellipsoid that results from the intersection of the plane  $h_1(\mathbf{m}) = \sigma_{\text{cubic}}^2 \chi_{31}^2(\alpha)$  with the surface  $h_2(\mathbf{m}) = \|\mathbf{d} - \mathbf{z}(\mathbf{m})\|^2$ . By projecting this region onto each of the axes, we obtain conservative  $1 - \alpha$  confidence intervals for each parameter. The right panel in Figure 6.1 shows the 95% projected intervals for  $g$  and  $C$ . The 95% and 99% projected



**Fig. 6.1** 95% confidence region that results from the intersection of the plane  $h_1(\mathbf{m}) = \sigma_{\text{cubic}}^2 \chi_{31}^2(0.05)$  with the surface  $h_2(\mathbf{m}) = \|\mathbf{d} - \mathbf{z}(\mathbf{m})\|^2$  for the two-parameter case.



**Table 6.1** Interval length and coverage results for the projected 95% intervals of  $g$  for the two- and three-parameter cases and likelihood (4.6).

$m$	$\hat{g} = \hat{g}_{ML}$		$\hat{g} = g_{ref}$	
	Coverage	Average length	Coverage	Average length
$(g, C)$	$97.02 \pm 0.01$	$0.66 \pm 0.01$	$96.92 \pm 0.01$	$0.68 \pm 0.01$
$(g, C, t_0)$	$97.72 \pm 0.01$	$2.60 \pm 0.04$	$97.52 \pm 0.01$	$2.60 \pm 0.04$

confidence intervals for  $g$  are

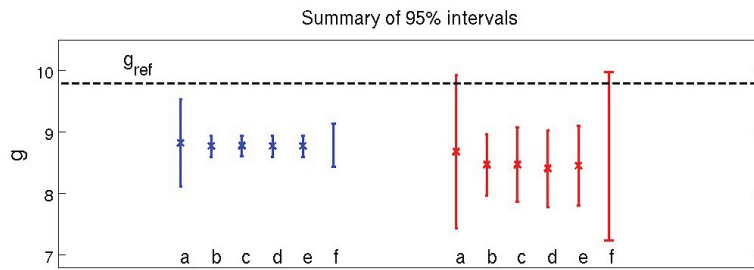
$$(6.1a) \quad \begin{array}{ccc} m & 95\% & 99\% \\ (g, C) & : [8.43, 9.14] & [8.37, 9.21] \end{array}$$

$$(6.1b) \quad \begin{array}{ccc} m & 95\% & 99\% \\ (g, C, t_0) & : [7.24, 9.98] & [7.04, 10.31]. \end{array}$$

Table 6.1 shows the corresponding coverage results again obtained using again a modification of Algorithm 1. As expected, these intervals are more conservative and wider than all the others. However, an advantage of this method is that we only need to use the forward model and find a rejection region defined by the noise distribution, and so there is no need to solve the inverse problem. Note also that one could use simulations to tune the value of  $\alpha$  so that the confidence level of the projected intervals is closer to 95%.

**7. Discussion.** The authors of A13 used real experimental data to obtain estimates of the value of  $g$  at the location where their experiment took place, College Station, TX. It was reassuring that the credible intervals under the three-parameter prior in A13 actually included  $g_{ref}$ . However, our validation analysis revealed that the spread of the posterior distribution, which allowed the inclusion of  $g_{ref}$  in the intervals, was actually caused by an artificially wide likelihood produced by an overcautious estimate of the measurement noise level. This is a reminder of the importance of checking the choice of likelihood function in addition to that of prior distributions. Once a better likelihood was selected, the credible intervals no longer covered  $g_{ref}$ . The prior distributions were chosen as is often done in practice; inequality constraints were enforced using flat priors. Given such ad-hoc choices, it is important to make sure the results are driven by the data and not by artificial information introduced by the priors. Prior and posterior predictive checks with the new likelihood did not reveal significant problems. In addition, the frequentist coverage of the credible intervals was near the target level. Thus, the priors do not seem to dominate the inference.

To validate the results, we also constructed confidence intervals using ML, non-linear LS, and test inversion. Figure 7.1 provides a summary of all the 95% intervals for  $g$  found in the previous sections. Once the new likelihood function is used, the Bayesian and non-Bayesian approaches give consistent results as judged by their frequentist behavior and even by the similarity of the actual intervals; excluding the intervals (6.1), which are very conservative, the other intervals are quite similar. Although not shown here, we also obtained similar results using a uniform or a Gaussian distribution with both different and equal  $\delta_i$ , in all cases modeled as symmetric about zero with a standard deviation that matches that of the actual noise in the data. Such agreement is expected in simple problems with only a few parameters to be estimated. It is in line with the Bernstein–von Mises theorem, which roughly states that if the amount of data (here  $n = 31$ ) is large compared to the number of unknown parameters (here  $p = 3$ ), then Bayesians and frequentists agree (see, for example,



**Fig. 7.1** 95% intervals for  $g$  for the case of two- (left group) and three-parameter (right group) vector  $\mathbf{m}$ : (a) credible intervals in A13; (b) and (c) are credible intervals (4.5) and (4.8), respectively; (d) ML intervals (5.5); (e) nonlinear LS intervals (5.8); (f) projected intervals (6.1).

[9, 11]). This is not necessarily true for more complex large-scale inverse problems [12, 10].

In addition to using frequentist methods to explore the validity of Bayesian procedures, we have seen how Bayesian methods can be used to derive procedures that have good frequentist properties. That is, just as ML provides a recipe to find estimators, we can often use Bayes theorem to do the same. Thus, Bayesian and frequentist methods of inference can be used in complementary ways.

It may be disappointing that our analysis here showed that the credible intervals are significantly smaller than those shown in A13 (primarily due to the fact that the measurement uncertainties are much smaller than assumed there) and do not include  $g_{\text{ref}}$ . However, this can be viewed in different ways. First, given that we know  $g_{\text{ref}}$ , this could indicate the presence of a systematic bias in the experiment that would not have been detectable using the original analysis without validation. Second, in practice parameters are of course unknown and the validation process shown here points out the importance of conducting a careful systematic error analysis to determine potential bias (and possible corrections). If the experiment described in A13 had been the first ever attempt to estimate  $g$  at College Station, the results would have underestimated its value, but we would have been unaware of this until a new, better experiment came along with a more accurate estimate. This is not unusual in science. Uncertainty estimates provide some guidance, but replication of the results with independent, better experiments is what scientists expect in order to accept a new result.

#### REFERENCES

- [1] M. ALLMARAS, W. BANGERTH, J. M. LINHART, J. POLANCO, F. WANG, K. WANG, J. WEBSTER, AND S. ZEDLER, *Estimating parameters in physical models through Bayesian inversion: A complete example*, SIAM Rev., 55 (2013), pp. 149–167.
- [2] L. D. BROWN, T. T. CAI, AND A. DASGUPTA, *Interval estimation for a binomial proportion*, Statist. Sci., 16 (2001), pp. 101–133.
- [3] B. P. CARLIN AND T. A. LOUIS, *Bayesian Methods for Data Analysis*, 3rd ed., CRC Press, Boca Raton, FL, 2008.
- [4] G. CASELLA AND R. L. BERGER, *Statistical Inference*, 2nd ed., Duxbury, Pacific Grove, CA, 2002.
- [5] A. C. DAVISON AND D. V. HINKLEY, *Bootstrap Methods and Their Applications*, Cambridge University Press, Cambridge, UK, 1997.
- [6] N. R. DRAPER AND H. SMITH, *Applied Regression Analysis*, 3rd ed., Wiley, New York, 1998.

- [7] B. EFRON AND D. V. HINKLEY, *Assessing the accuracy of the maximum likelihood estimator: Observed versus expected Fisher information*, *Biometrika*, 65 (1978), pp. 457–482.
- [8] B. EFRON AND R. TIBSHIRANI, *An Introduction to the Bootstrap*, Chapman & Hall, London, 1993.
- [9] T. S. FERGUSON, *A Course in Large Sample Theory*, Chapman & Hall, London, 1996.
- [10] C. FOX, Y. MARZOUK, AND L. TENORIO, *An Introduction to Data Analysis and Uncertainty Quantification for Inverse Problems*, in preparation.
- [11] D. FREEDMAN, *Some issues in the foundation of statistics*, *Found. Sci.*, 1 (1995), pp. 19–39.
- [12] D. FREEDMAN, *On the Bernstein-von Mises theorem with infinite-dimensional parameters*, *Ann. Statist.*, 27 (1999), pp. 1119–1140.
- [13] A. GELMAN, J. CARLIN, H. S. STERN, AND D. B. RUBIN, *Bayesian Data Analysis*, 2nd ed., Chapman & Hall, London, 2003.
- [14] R. A. JOHNSON AND D. W. WICHERN, *Applied Multivariate Analysis*, 6th ed., Prentice-Hall, Englewood Cliffs, NJ, 2007.
- [15] E. L. LEHMANN AND J. P. ROMANO, *Testing Statistical Hypotheses*, Springer, New York, 2005.
- [16] J. RICE, *Mathematical Statistics and Data Analysis*, 3rd ed., Duxbury, Belmont, CA, 2007.
- [17] F. J. SAMANIEGO, *A Comparison of the Bayesian and Frequentist Approaches to Estimation*, Springer, New York, 2010.
- [18] G. A. F. SEBER AND C. J. WILD, *Nonlinear Regression*, Wiley, New York, 2003.
- [19] P. B. STARK, *Inference in infinite-dimensional inverse problems: Discretization and duality*, *J. Geophys. Res.*, 97 (1992), pp. 14055–14082.
- [20] P. B. STARK AND L. TENORIO, *A primer of frequentist and Bayesian inference in Large-Scale Inverse Problems and Quantification of Uncertainty*, Wiley, 2011, pp. 9–32.
- [21] L. TENORIO, F. ANDERSSON, M. DE HOOP, AND P. MA, *Data analysis tools for uncertainty quantification of inverse problems*, *Inverse Problems*, 27 (2011), 045001.