



ISTITUTO NAZIONALE DI RICERCA METROLOGICA Repository Istituzionale

Bayesian model selection applied to linear regressions with weighted data

This is the author's accepted version of the contribution published as:

Original

Bayesian model selection applied to linear regressions with weighted data / Mana, Giovanni; Massa, Enrico; Sasso, Carlo Paolo. - In: METROLOGIA. - ISSN 0026-1394. - 56:2(2019). [10.1088/1681-7575/ab0338]

Availability:

This version is available at: 11696/60147 since: 2020-12-21T19:24:51Z

Publisher:

IOP

Published

DOI:10.1088/1681-7575/ab0338

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

Institute of Physics Publishing Ltd (IOP)

IOP Publishing Ltd is not responsible for any errors or omissions in this version of the manuscript or any version derived from it. The Version of Record is available online at DOI indicated above

(Article begins on next page)

Bayesian model selection applied to linear regressions with weighted data

G Mana, E Massa and C P Sasso

INRIM - Istituto Nazionale di Ricerca Metrologica, Str. delle Cacce 91, 10135
Torino, Italy

E-mail: g.mana@inrim.it

Abstract. This paper uses the Bayesian model-selection to find the linear regression that is most supported by the data. It uses a hierarchical model to improve and to extend to weighted data a previous investigation. As an application example, it shows how the results can be used to investigate the consistency of measurements carried out to determine the central second-moment of the angular power-spectrum of a laser beam.

Submitted to: *Metrologia*

PACS numbers: 02.50.Cw, 02.50.Tt, 07.05.Kf

1. Introduction

A problem in linear regression is to test the goodness of fit and to select the model most supported by the data. If the data explanation is uncertain, for instance, when building calibration curves, the problem is to choose how many basis functions to include in the regression. If the data uncertainties are underestimated, different hypotheses are possible, for instance, in the presence of random or fixed effects [1, 2, 3].

To make meaningful assessments of the odds of competing explanations one must assign a prior distribution to the regression parameters in such a way that the posterior probabilities of regressions that differ only by their parametrisation are the same. This requirement was not considered in previous works [3, 4, 5, 6]. These prior distributions (named after H. Jeffreys) are proportional to the invariant volume-element of the model manifold [7, 8], where the regression parameters are the manifold coordinates and whose metric is the Fisher information [9, 10, 11, 12].

A difficulty impedes to pursue this line of thought: Jeffrey's distribution of the regression parameters is improper. To go around this problem, we build on previous proposals [4, 6] and use a hierarchical model to improve and extend to weighted and correlated data the results given in [6].

Section 2 frames the study and introduces the representations of the data and design-matrix that make the algebra the simplest. The Bayesian procedure to compare competing models is shortly summarised in section 3 and applied to the problem at hand in section 4. Eventually, a numerical example and an application to experimental data demonstrate the practical use of model selection in data analysis.

2. Linear regressions

Let us explain N observations, $\mathbf{y} = \{y_1, y_2, \dots, y_N\}^T$, by a set of competing (mutually exclusive) linear models

$$\mathbf{y} = W_l \mathbf{b} + \mathbf{u}, \quad (1)$$

where l is the regression order, \mathbf{u} are zero-mean errors having variance-covariance matrix Σ_u , $W_l = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l\}$ are full rank $N \times l$ design matrices, \mathbf{w}_i are the model basis-vectors, and $\mathbf{b} = \{b_1, b_2, \dots, b_l\}^T$ are model parameters. The problem is to find the model most supported by the data. The basis functions, whose discrete samples \mathbf{w}_i form the basis vectors, may be polynomials but, in general, they are any set of linearly independent functions. For the sake of simplicity, we indexed the design matrices W_l by the regression order, but competing regressions having the same order are not excluded.

In order to avoid awkward algebra, it is convenient to normalize the data in such a way that they have unit variance and are uncorrelated. This corresponds to use an orthonormal basis in the \mathbf{y} space and it is done by the transformations $\mathbf{h} = U^{-1}\mathbf{y}$, $V_l = U^{-1}W_l$, and $\boldsymbol{\epsilon} = U^{-1}\mathbf{u}$, where the lower triangular matrix U is the Cholesky factor of Σ_u , that is, $\Sigma_u = UU^T$ [13]. Hence, we rewrite the data model (1) as $\mathbf{h} = V_l \mathbf{b} + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon}$ are zero-mean uncorrelated errors having unit variance-covariance matrix, $V_l = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l\}$, and \mathbf{v}_i are the new design matrices and basis vectors.

Since we are not interested in any specific parametrisation, a further simplification is obtained by using an orthonormal basis in the subspace spanned by $V_l \mathbf{b}$, which is done by orthonormalizing the $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l\}$ set and by scaling the parameters according to the decomposition $V_l = Q_l R$, where Q_l is column-orthonormal, $Q_l^T Q_l = \mathbb{1}$, $\mathbb{1}$ is the unit matrix, and R is upper triangular [13].

Eventually, the maximum-entropy distribution of the normalized data, constrained by $\mathbf{h} = Q_l \mathbf{a} + \boldsymbol{\epsilon}$ and $\Sigma_h = \mathbb{1}$ is

$$L(\mathbf{h}|\mathbf{a}, l) = \sqrt{\frac{1}{(2\pi)^N}} \exp \left[-\frac{|\mathbf{h} - Q_l \mathbf{a}|^2}{2} \right], \quad (2)$$

where $\mathbf{a} = R\mathbf{b} = \{a_1, a_2, \dots, a_l\}^T$. Since we are not interested in a specific model-parametrisation, but only in the rank l of the design matrix, no generality is lost.

3. Bayesian model comparison

Let $\{Q_l\}$ be a complete set of mutually exclusive models – indexed by the regression order l – that compete to explain the dataset \mathbf{h} . Prior the measurement, the joint distribution of the data, model, and model parameters \mathbf{a} can be written in terms of conditional distributions as

$$P(\mathbf{h}, \mathbf{a}, l) = L(\mathbf{h}|\mathbf{a}, l) \pi(\mathbf{a}|l) \Pi(l), \quad (3)$$

where the likelihood $L(\mathbf{h}|\mathbf{a}, l)$ is the sampling distribution of \mathbf{h} , $\pi(\mathbf{a}|l)$ is the pre-data distribution of the model parameters, and $\Pi(l)$ is the prior probability of Q_l being true. Accordingly, l – hence, the model – is sampled from $\Pi(l)$; then, the model parameters are sampled from $\pi(\mathbf{a}|l)$; eventually, the data are sampled from $L(\mathbf{h}|\mathbf{a}, l)$. The prior distribution $\Pi(l)$ must synthesize information available about the data model. In the absence of information, the probabilities that maximizes the $\Pi(l)$ entropy are equal.

By marginalisation of $P(\mathbf{h}, \mathbf{a}, l)$ over the model parameters and conditioning on the data, the posterior probability of the l -th model provided by the data is

$$P(l|\mathbf{h}) = \frac{Z(\mathbf{h}|l)\Pi(l)}{\sum_l Z(\mathbf{h}|l)\Pi(l)}. \quad (4)$$

Here, the marginal likelihood (or evidence),

$$Z(\mathbf{h}|l) = \int_A L(\mathbf{h}|\mathbf{a}, l)\pi(\mathbf{a}|l)d\mathbf{a}, \quad (5)$$

where the integration is carried out on the parameter space, is the sampling distribution of the data given the model, independently of what the parameter values may be. It summarises the uncertainty about what is the model explaining the data.

Not unlike the case where the Bayes theorem is applied to determine the posterior probabilities of the model-parameter values, the distribution (4) embeds all the statements about the uncertainty after the data have been observed. In parameter estimation, uncertainty summaries are the posterior standard deviation and credible intervals. Also, the optimal estimate under quadratic loss of a parameter function (i.e., a measurand) is its expectation with respect the posterior distribution of the parameter values. The same is true for the posterior model-distribution. A simple way to select a model is to chose the mode. When $P(l|\mathbf{h})$ does not allow for a reliable comparison of posterior probabilities and no single model stands out, we can consider all the models and average over them. Whether and how to synthesize $P(l|\mathbf{h})$ – e.g., by selection or averaging – depends on decision theory considerations.

To investigate how much $P(l|\mathbf{h})$ depends on the assumptions made – that is, on the priors $\Pi(l)$ and $\pi(\mathbf{a}|l)$ and likelihood $L(\mathbf{h}|\mathbf{a}, l)$ – we can apply the Bayesian analysis to all plausible combinations of priors and likelihood. Then, the $P(l|\mathbf{h})$ variability with these changes delivers information about how much confidence we can place on the model W_l . This sensitivity analysis is out of the scope of this manuscript. Besides, the $\Pi(l)$ and $L(\mathbf{h}|\mathbf{a}, l)$ distributions have been uniquely set by the available information via entropy maximization. Furthermore, sections 4.1 and 4.2, will prove that also $\pi(\mathbf{a}|l)$ is uniquely imposed to result in equal posterior probabilities of models that differ only in parametrization.

The marginal likelihood and posterior model-probability are proportional to the ability of the model to explain the data – the higher the fitness, the higher $Z(\mathbf{h}|l)$ – but inversely proportional to the volume of the parameter space – the higher the model freedom, the lesser $Z(\mathbf{h}|l)$. Therefore, models with fewer parameters, which have a smaller parameter space, are preferred. This bias is known as the Ockham's razor and penalises the models having greater freedom in explaining the data.

4. Application to linear regression

4.1. Prior distribution of the model parameters

In order to find the model probabilities via (4) and (5), we need the pre-data distribution $\pi(\mathbf{a}|l)$. The probability calculus updates prior distributions into posterior ones, but it does not tell how to assign the prior probabilities.

In the absence of pre-data information, the posterior model-probabilities $P(l|\mathbf{h})$ and, consequently, the marginal likelihood $Z(\mathbf{h}|l)$ must be independent of the model parameters, that is, they must be invariant for model reparametrisations and transform as probability distributions. For instance, if the explaining model is a polynomial, (4)

and (5) must depend only on the polynomial degree l and not on what basis, orthogonal or not, is used.

This is ensured by the Jeffreys' prior $\pi_J(\mathbf{a}|l) \propto \sqrt{|F(\mathbf{a}|l)|}$, where $F(\mathbf{a}|l)$ is the Fisher information that the data carry about the parameters and the vertical bars denote the determinant. To calculate $\pi_J(\mathbf{a}|l)$, let us observe that the log-likelihood Jacobian of (2) is [14]

$$J = \partial_{\mathbf{a}} [\ln(L)] = -\partial_{\mathbf{a}} (\mathbf{x}^T \mathbf{x}) / 2 = -\mathbf{x}^T (\partial_{\mathbf{a}} \mathbf{x}) = \mathbf{x}^T Q_l \quad (6)$$

where $\mathbf{x} = \mathbf{h} - Q_l \mathbf{a}$. Hence,

$$F(\mathbf{a}|l) = -\langle \partial_{\mathbf{a}} J \rangle_{h|a,l} = Q_l^T Q_l = \mathbb{1} \quad (7)$$

where $\langle \rangle_{h|a,l}$ indicates the average with respect to the data conditioned to the model and model parameters. Since (7) does not depend on the model parameters, $\pi_J(\mathbf{a}|l) \propto \text{const.}$ is improper.

The use of an improper prior to calculate the posterior distribution of the model parameters,

$$P(\mathbf{a}|\mathbf{h}, l) = \frac{L(\mathbf{h}|\mathbf{a}, l) \pi_J(\mathbf{a}|l)}{Z(\mathbf{h}|l)}, \quad (8)$$

is justified by showing that $P(\mathbf{a}|\mathbf{h}, l)$ is the limit of the posteriors obtained from proper priors defined on increasingly large bounded-supports [15]. However, this argument encodes that $|\mathbf{a}|$ is greater than any positive number. This is irrelevant to the \mathbf{a} posterior, but, it is not so for $Z(\mathbf{h}|l)$ and $P(l|\mathbf{h})$. In fact, since $Z(\mathbf{h}|l)$ and $P(l|\mathbf{h})$ are defined only up to model-dependent, but unknown, scale factors, an improper prior makes their calculation meaningless.

4.2. Hierarchical model

A way to escape from this difficulty is to complement the problem with measurable prior information, which information makes the $Z(\mathbf{h}|l)$ and $P(l|\mathbf{h})$ invariance irrelevant. Since we know in advance that both the expected value and variance of the data are bounded [4], we can encode in the parameter prior a prejudice towards repeated measurements of the same quantity, whose value, α , is unknown. This prejudice is modelled by setting $\mathbf{w}_1 = \{1, 1, \dots, 1\}^T$ in (1), where \mathbf{w}_1 is the first column of W_l , and $\langle \mathbf{y} \rangle_{b,u} = \alpha \mathbf{w}_1$. The subscript indicates that we take the mean over the joint (prior) distribution of \mathbf{b} and \mathbf{u} .

After switching to the normalized data and orthonormal basis, the relationships $|\mathbf{q}_1| = 1$ and $Q_l Q_l^T \mathbf{q}_1 = \mathbf{q}_1$ hold. It is convenient to ensure that the estimate $\hat{a}_1 = \mathbf{q}_1^T \mathbf{h}$ is proportional to the weighted mean of the data,

$$y_{LS} = (\mathbf{w}_1^T \Sigma_u^{-1} \mathbf{w}_1)^{-1} \mathbf{w}_1^T \Sigma_u^{-1} \mathbf{y}. \quad (9a)$$

This is achieved by starting the orthonormalisation from $\mathbf{v}_1 = U^{-1} \mathbf{w}_1$; hence, $\mathbf{q}_1 = \mathbf{v}_1 / |\mathbf{v}_1|$. In fact,

$$y_{LS} = (\mathbf{v}_1^T \mathbf{v}_1)^{-1} \mathbf{v}_1^T \mathbf{h} = \hat{a}_1 / |\mathbf{v}_1|, \quad (9b)$$

where we used $\Sigma_u^{-1} = U^{-T} U^{-1}$, $U^{-1} \mathbf{y} = \mathbf{h}$, and $\mathbf{q}_1^T \mathbf{q}_1 = 1$.

That done, the repeated measurements of a constant prejudice is still encoded by $\langle \mathbf{h} \rangle_{a,\epsilon} = \alpha \mathbf{q}_1$, where we redefined the hyperparameter α . Hence, $Q \langle \mathbf{a} \rangle_a = \langle \mathbf{h} \rangle_{a,\epsilon} = \alpha \mathbf{q}_1$ and

$$\langle \mathbf{a} \rangle_a = \alpha Q_l^T \mathbf{q}_1 = \alpha \mathbf{p}_1, \quad (10)$$

where $\mathbf{p}_1 = Q_l^T \mathbf{q}_1 = \{1, 0, \dots, 0\}^T$ and we used $Q_l^T Q_l = \mathbb{1}$. Next, we encoded the bounded data-variance prejudice by

$$\Sigma_a = \langle \mathbf{a} \mathbf{a}^T \rangle_a = (\beta^2 - 1) \mathbb{1}, \quad (11)$$

where we used $\langle \mathbf{h} \mathbf{h}^T \rangle_{a, \epsilon} = \beta^2 \mathbb{1} = Q_l \Sigma_a Q_l^T + \mathbb{1}$ and $\beta > 1$.

The prior distribution of the model parameters constrained by (10) and (11) and maximizing the relative entropy with respect to the uniform distribution is

$$\pi(\mathbf{a} | \alpha, \beta, l) = \sqrt{\frac{1}{(2\pi)^l (\beta^2 - 1)^l}} \exp \left[-\frac{|\mathbf{a} - \alpha \mathbf{p}_1|^2}{2(\beta^2 - 1)} \right], \quad (12)$$

where the hyper-parameters α and β are unknown. It is worth noting that, when $\beta \rightarrow \infty$, $\pi(\mathbf{a} | \alpha, \beta, l)$ reproduces the uniform prior.

4.3. Sampling distribution of the data

The distribution of the normalized data, given the hyper-parameters α and β and the orthonormal model Q_l ,

$$\begin{aligned} \mathcal{L}(\mathbf{h} | \alpha, \beta, l) &= \int_{\mathbf{R}^l} L(\mathbf{h} | \mathbf{a}, l) \pi(\mathbf{a} | \alpha, \beta, l) d\mathbf{a} \\ &= \sqrt{\frac{1}{(2\pi)^n \beta^{2l}}} \exp \left(-\frac{(\mathbf{h} - \alpha \mathbf{q}_1)^T [\mathbb{1} + (\beta^2 - 1) Q_l Q_l^T]^{-1} (\mathbf{h} - \alpha \mathbf{q}_1)}{2} \right), \end{aligned} \quad (13)$$

is obtained by marginalization over the model parameters. The relevant integration is carried out in the Appendix A.

4.4. Prior distribution of the hyper-parameters

To ensure that $P(l | \mathbf{h})$ is invariant for model reparametrisations, we must compute the Jeffreys' prior $\pi_J(\alpha, \beta | l)$. As it is shown in the Appendix B, the Fisher information that \mathbf{h} carries about α and β is

$$F(\alpha, \beta) = \frac{1}{\beta^2} \begin{pmatrix} 1 & 0 \\ 0 & 2l \end{pmatrix}. \quad (14)$$

Therefore, the Jeffrey's distribution of α and β given the Q_l model is

$$\pi_J(\alpha, \beta | l) \propto \sqrt{|F(\alpha, \beta | l)|} \propto 1/\beta^2, \quad (15)$$

whose support is $\beta > 1$, $-\infty < \alpha < +\infty$. Although it is improper, since now the same normalizing-factor is included in all the marginal likelihoods, (15) does not impede the Q_l comparison. In fact, after the normalization, the model probabilities (4) do not depend on the $\pi_J(\alpha, \beta | l)$ support in the α space.

4.5. Marginal likelihood

The computation of the marginal likelihood (the data sampling-distribution given W_l) is given in Appendix C. Although it was derived using uncorrelated data having the same unit variance and an orthonormal basis in the model subspace, the result,

$$Z(\mathbf{y} | l) \propto \frac{e^{-|\hat{\mathbf{u}}|^2/2} \gamma(l/2, |\hat{\mathbf{y}}'|^2/2)}{(|\hat{\mathbf{y}}'|/\sqrt{2})^l}, \quad (16)$$

is independent of the transformations made and it is written terms of the weighted least-squares estimate of the original data

$$\hat{\mathbf{y}}' = W_l(W_l^T \Sigma_u^{-1} W_l)^{-1} W_l^T \Sigma_u^{-1} \mathbf{y}', \quad (17)$$

where $\mathbf{y}' = \mathbf{y} - \mathbf{w}_1 y_{LS}$ and y_{LS} is the weighted mean (9a), and the residuals $\hat{\mathbf{u}}$. Unessential factors independent of the regression order l and data have been left out and $\gamma(s, z_2)$ is the lower incomplete gamma function [16]. The weighted norms $|\hat{\mathbf{u}}|^2$ and $|\hat{\mathbf{y}}'|^2$ are calculated according to the Σ_u^{-1} metric, that is, $|\hat{\mathbf{u}}|^2 = \hat{\mathbf{u}}^T \Sigma_u^{-1} \hat{\mathbf{u}}$ and $|\hat{\mathbf{y}}'|^2 = \hat{\mathbf{y}}'^T \Sigma_u^{-1} \hat{\mathbf{y}}'$.

Eventually, with the assumed uniform $\Pi(l)$ prior distribution, the posterior probability of each W_l model is obtained by the normalization to $\sum_l Z(\mathbf{y}|l) = 1$.

When we consider independent and identically distributed data $\Sigma_u = \sigma_0^2 \mathbf{1}$ and

$$Z(\mathbf{y}|l) \propto \frac{e^{-|\hat{\mathbf{u}}|^2/(2\sigma_0^2)} \gamma[l/2, |\hat{\mathbf{y}}'|^2/(2\sigma_0^2)]}{[|\hat{\mathbf{y}}'|/(\sqrt{2}\sigma_0)]^l}, \quad (18)$$

where, now, y_{LS} is the arithmetic mean of the data, $\hat{\mathbf{y}}' = W_l(W_l^T W_l)^{-1} W_l^T \mathbf{y}'$ is the ordinary least-squares estimate, and $|\hat{\mathbf{u}}|^2 = \hat{\mathbf{u}}^T \hat{\mathbf{u}}$ and $|\hat{\mathbf{y}}'|^2 = \hat{\mathbf{y}}'^T \hat{\mathbf{y}}'$ are the norms of the residuals and data-estimates calculated according the Euclidean metric. We observe that (18) is the same as the equation (B.4) in [6]. In fact, since $\hat{\mathbf{y}}^T(\mathbf{y} - \hat{\mathbf{y}}) = 0$ (because $\hat{\mathbf{y}}$ is the \mathbf{y} projection in the $W_l \mathbf{b}$ sub-space), $\mathbf{y}^T \hat{\mathbf{u}} = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) = |\hat{\mathbf{u}}|^2$.

Extending the comparison to non-linear models would require determining the Jeffreys' prior of the model parameters, which, in general, is a difficult problem to solve and devoid of a general solution. However, since it does not depend on the W_l representation, but only on the regression order, residuals, and central moments of the data, (16) can be heuristically applied also to non-linear models, provided that their curvature in the parameter subspace of interest can be neglected – which means that a linear approximation holds.

4.6. Asymptotic behaviour

For a large data sample, provided $|\hat{\mathbf{y}}'|^2 \gg l$,

$$\ln [Z(\mathbf{y}|l)] \propto -|\hat{\mathbf{u}}|^2/2 - l \ln(|\hat{\mathbf{y}}'|), \quad (19a)$$

where we omitted the terms independent of l and used $\gamma(l/2, |\hat{\mathbf{y}}'|^2/2) \approx \Gamma(l/2)$ and $\Gamma(l/2) \ll |\hat{\mathbf{y}}'|^l$ [17]. Furthermore, for large $l \ll |\hat{\mathbf{y}}'|^2$, $|\hat{\mathbf{y}}'|^2$ is independent of the regression order and proportional to the sample size N . Hence,

$$\ln [Z(\mathbf{y}|l)] \propto -|\hat{\mathbf{u}}|^2/2 - l \ln(N)/2, \quad (19b)$$

always omitting the terms independent of l .

The model having the maximum posterior-probability maximises the log-likelihood. The asymptotic behaviour (19a) shows that the optimal order minimises the weighted residuals by keeping the number of free parameters as small as possible, thus handicapping the overfitting. Among the regressions having the same order, the most supported is that whose associated sum of the (weighted) squared residuals is minimum.

To compare (16) with the classical indicators, we compare (19a-b) with the marginal likelihoods of the asymptotic information criteria due to Akaike [18, 19],

$$\ln [Z_{AIC}(\mathbf{y}|l)] \propto -|\hat{\mathbf{u}}|^2/2 - l, \quad (20a)$$

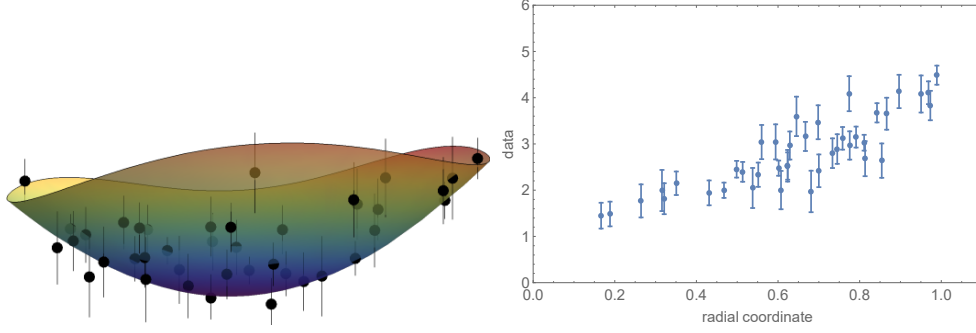


Figure 1. Left: three dimensional scatter plot of the data sampled from (21). The bars indicate 95% confidence intervals. The best-fit regression to the data is also shown. Right: radial plot of the same data.

and Schwarz [20, 21, 22],

$$\ln [Z_{BIC}(\mathbf{y}|l)] \propto -|\hat{\mathbf{u}}|^2/2 - l \ln(N)/2, \quad (20b)$$

where we omitted the terms that are independent of the regression order (which does not matter, because all the models are Gaussian). The identity of (19b) and (20b) shows that the model probabilities derived from (16) shed light and encompass these (approximate) criteria.

5. Numerical example

Figure 1 shows a set of 40 data randomly sampled in the unit disk from the aberrated wavefront

$$y = b_1 Z_0^0(x_1, x_2) + b_2 Z_1^1(x_1, x_2) + b_3 Z_2^0(x_1, x_2) + b_4 Z_3^{-1}(x_1, x_2), \quad (21)$$

where $Z_n^m(x_1, x_2)$ are Zernike polynomials, $\{Z_0^0, Z_1^1, Z_2^0, Z_3^{-3}\}$ is the basis set, and the parameter values are $\mathbf{b} = \{3, 0.5, 1.5, 0.5\}$. Zero mean uncorrelated Gaussian errors – having random standard-deviations $\sigma_0(i)$ ranging from 5% to 23% – were added to the data.

To explain the data, we carried out regressions (having order ranging from $l = 1$ to $l = 11$) using the 1024 basis sets corresponding to subsets – containing Z_0^0 – of

$$\{Z_0^0, Z_1^{-1}, Z_1^1, Z_2^0, Z_3^{-3}, Z_3^3, Z_3^{-1}, Z_3^1, Z_4^0\}. \quad (22)$$

We calculated the marginal likelihood and posterior probability of each basis by application of (16) and (4), where the bases were assumed to have the same prior probability.

To assess the efficiency of (16) in detecting the right model, we repeated the regressions 100 times – by using 100 independent datasets – and averaged the posterior model-probabilities. After sorting the results into decreasing probabilities, Fig. 2 (left) shows the top eight posterior-probabilities. The most sustained basis is that used to generate the data.

In a second test, we used exponentially correlated errors. Therefore, the variance-covariance matrix of the data is

$$\Sigma_u(i, j) = \sigma_0(i)\sigma_0(j)e^{-5r_{ij}}, \quad (23)$$

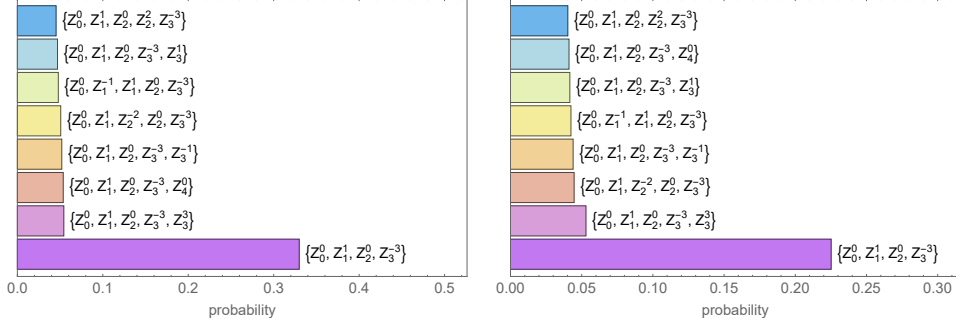


Figure 2. Top eight probabilities of the subsets of (22) averaged over 100 Monte Carlo simulations. Left: uncorrelated data. Right: exponentially correlated data.

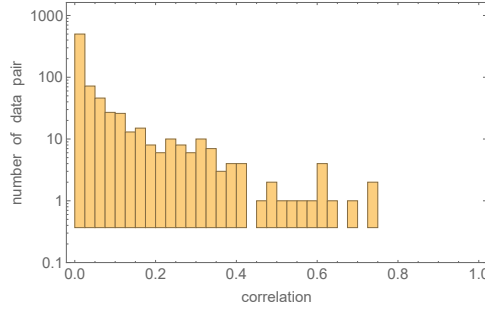


Figure 3. Histogram of the correlations between the data sampled from (21).

where r_{ij} is the distance of the datum i from the j one. Figure 3 shows an example of the data correlations. Also in this case, we repeated the regressions 100 times – by using 100 independent datasets – and averaged the posterior model-probabilities. Figure 2 (right) shows the result. Though the mean probability of the $\{Z_0^0, Z_1^1, Z_2^0, Z_3^{-3}\}$ basis decreases from 33% to 23%, it is still the most fostered by the data.

6. Application example

The marginal likelihood given in (16) was used to investigate doubtful measurements carried out to determine the central second-moment of the angular power-spectrum of a laser beam, which is a pivotal quantity when correcting for the diffraction effects the length measurements carried out by optical interferometry. We do not want here describe these measurements nor draw conclusions about them, but only show how the results obtained can be used in practice.

The second moment is measured by imaging the power spectrum in the focal plane of a telescope and by numerically integrating the normalised image by a sum over the camera pixels [23]. We repeated the measurement by placing the telescope at 137 cm, 394 cm, and 1047 cm distance from the beam source and grabbed defocused images by putting the camera in the focal plane and at $(-15, -10, -5, +5, +10, +15)$ mm from it. The results are shown in Fig. 4. The paraxial propagation of a scalar beam in free space predicts that the focal-plane images do not depend on the telescope

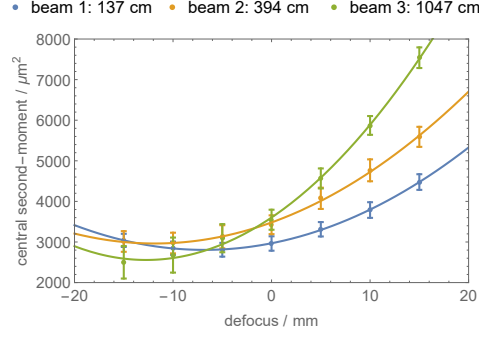


Figure 4. Central second-moment of the defocused images. Each set of seven data corresponds to a different telescope distance. The solid lines are the best parabolas fitting the data.

distance and that the seven central-moments depend quadratically on the defocus.

Figure 4 shows that the best-fit parabolas to the data do not intersect where the defocus is null, as we would have expected. To quantify how much the data support this observation and the quadratic dependence on the defocus, we compared three hypotheses: the data belong to i) non-intersecting parabolas, that is,

$$\mathbf{y} = W_9 \mathbf{b} + \mathbf{u} \quad (24a)$$

or

$$y_{ij} = b_{i0} + b_{i1}z_j + b_{i2}z_j^2 + u_{ij}, \quad (24b)$$

where the data and parameters are assembled in \mathbf{y} and \mathbf{b} , z is the defocus, $i = 1, 2, 3$ labels the data series, and $j = 1, 2, \dots, 7$ labels the defocus; ii) non-intersecting quartics functions, that is,

$$\mathbf{y} = W_{12} \mathbf{b} + \mathbf{u} \quad (25a)$$

or

$$y_{ij} = b_{i0} + b_{i1}z_j + b_{i2}z_j^2 + b_{i4}z_j^4 + u_{ij}; \quad (25b)$$

and iii) parabolas intersecting at null defocus, that is

$$\mathbf{y} = W_7 \mathbf{b} + \mathbf{u} \quad (26a)$$

or

$$y_{ij} = b_0 + b_{i1}z_j + b_{i2}z_j^2 + u_{ij}. \quad (26b)$$

The design matrices are given in the Appendix D.

With assumed equal prior probabilities of the models – after removing the weighted mean from the data and calculating their weighted least-square estimates and residuals – the marginal likelihood (16) and its normalization to $\sum_l Z(\mathbf{y}|l) = 1$ delivers the posterior probabilities shown in Fig. 5 (left). It confirms that the second moment depends quadratically on the defocus and the non-intersecting parabolas hypothesis overcomes by a 99.99 % probability the other two. We don't yet have an explanation for this unexpected behaviour.

For a Gaussian beam, the second moment and its curvature at the vertex are constrained by the unity value of the quality factor. Since a lens does not change the unit value of the quality factor, the beams downstream of the telescope lens must be

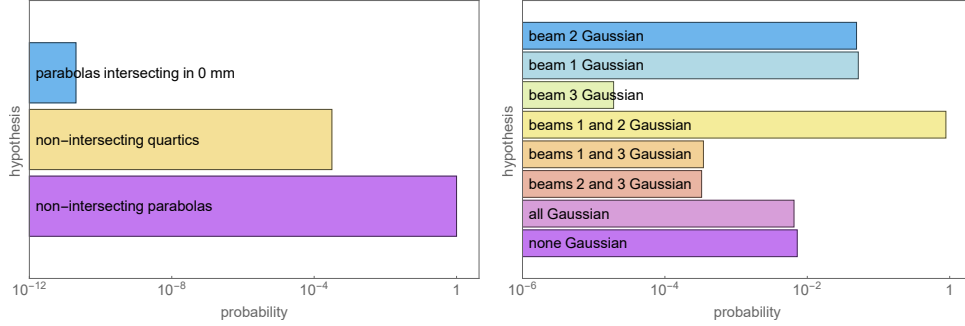


Figure 5. Posterior probabilities of the listed explanations of the data shown in Fig. 4

all Gaussian or all paraxial. Therefore, as an additional application of (16) (this time heuristically comparing non-linear models), we calculated the probabilities that none, one, two, or all the downstream beams are Gaussian.

A unit quality factor implies that the second moment propagates as $y^2 = a + c(z - z_0)^2$, where $k\sqrt{ac} = 1$ constrains the a and c parameters, k is the wavenumber, and z_0 the vertex position. Hence, the data models are

$$y_{ij} = a_i + \frac{(z_j - z_{i0})^2}{a_i k^2} + u_{ij}, \quad (27)$$

if the i -th beam is Gaussian, and

$$y_{ij} = a_i + c_i(z_j - z_{i0})^2 + u_{ij}, \quad (28)$$

if the i -th beam is paraxial. Eventually, eight (non-linear) data models were built by arranging in a eight systems three equations (corresponding to $i = 1, 2, 3$) drawn from (27) and/or (28).

Figure 5 (right) shows the posterior probabilities, where all the explanations were assumed to have the same prior probability. Also in this case (16) identifies a single explanation, assigns it an 88.4 % probability, and points out a severe inconsistency. This result was confirmed by calculating the M^2 value of the three beams.

7. Conclusions

This paper showed how to select the linear regression most supported by the data and extended previous works that did not account for weights and correlations [4, 6]. The selection requires that the marginal likelihood and posterior regression probabilities be independent of the parametrisation used, which previous investigations did not consider [3, 5]. To comply with this requirement, we derived Jeffreys priors from the volume element of each model-manifold equipped with the information metric. Also, to avoid inconsistencies due to improper priors, we used hierarchical modelling.

Although the computation of the marginal likelihood required complex integrations, the final formula (16), which can also be heuristically used for non-linear models, involves only standard least-squares estimates and simple algebra. Our result makes Bayesian model selection easily accessible and extends this toolbox to the analysis of arbitrarily correlated and weighted data.

Appendix A. Data likelihood

To calculate the sampling-distribution $\mathcal{L}(\mathbf{h}|\alpha, \beta, l)$ we observe that

$$|\mathbf{h} - Q_l \mathbf{a}|^2 = |\mathbf{a} - \hat{\mathbf{a}}|^2 + |\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2, \quad (\text{A.1})$$

where $\hat{\mathbf{a}} = Q_l^T \mathbf{h} = \{\hat{a}_1, \hat{a}_2, \dots, \hat{a}_l\}^T$ are the least squares estimates of the model parameters. Furthermore,

$$|\mathbf{a} - \alpha \mathbf{p}_1|^2 = |\mathbf{a}'|^2 + (a_1 - \alpha)^2. \quad (\text{A.2})$$

and

$$|\mathbf{a} - \hat{\mathbf{a}}|^2 = |\mathbf{a}' - \hat{\mathbf{a}}'|^2 + (a_1 - \hat{a}_1)^2, \quad (\text{A.3})$$

where $\hat{a}_1 = \mathbf{q}_1^T \mathbf{h}$, $\mathbf{a}' = \{a_2, a_3, \dots, a_l\}^T$ and $\hat{\mathbf{a}}' = \{\hat{a}_2, \hat{a}_3, \dots, \hat{a}_l\}^T$. Hence, by carrying out the needed integrations with the aid of Mathematica[®] [17],

$$\begin{aligned} \mathcal{L}(\mathbf{h}|\alpha, \beta, l) &= \sqrt{\frac{1}{(2\pi)^{n+l}(\beta^2 - 1)^l}} \int_{-\infty}^{+\infty} da_1 \int_{\mathbf{R}^{l-1}} \exp\left(-\frac{|\mathbf{h} - Q_l \mathbf{a}|^2}{2}\right) \exp\left[-\frac{|\mathbf{a} - \alpha \mathbf{p}_1|^2}{2(\beta^2 - 1)}\right] d\mathbf{a}' \\ &= \sqrt{\frac{1}{(2\pi)^{n+l}(\beta^2 - 1)^l}} \exp\left(-\frac{|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2}{2}\right) \int_{-\infty}^{+\infty} \exp\left(-\frac{(a_1 - \hat{a}_1)^2}{2}\right) \exp\left(-\frac{(a_1 - \alpha)^2}{2(\beta^2 - 1)}\right) da_1 \\ &\quad \times \int_{\mathbf{R}^{l-1}} \exp\left(-\frac{|\mathbf{a}' - \hat{\mathbf{a}}'|^2}{2}\right) \exp\left(-\frac{|\mathbf{a}'|^2}{2(\beta^2 - 1)}\right) d\mathbf{a}' \\ &= \sqrt{\frac{1}{(2\pi)^{(n+1)\beta^2(l-1)(\beta^2 - 1)}}} \exp\left(-\frac{|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2}{2} - \frac{|\hat{\mathbf{a}}'|^2}{2\beta^2}\right) \\ &\quad \times \int_{-\infty}^{+\infty} \exp\left(-\frac{(a_1 - \hat{a}_1)^2}{2}\right) \exp\left(-\frac{(a_1 - \alpha)^2}{2(\beta^2 - 1)}\right) da_1 \\ &= \sqrt{\frac{1}{(2\pi)^n \beta^{2l}}} \exp\left(-\frac{|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2}{2} - \frac{|\hat{\mathbf{a}}'|^2}{2\beta^2} + \frac{(\hat{a}_1 - \alpha)^2}{2\beta^2}\right) \\ &= \sqrt{\frac{1}{(2\pi)^n \beta^{2l}}} \exp\left(-\frac{|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2}{2} - \frac{|\hat{\mathbf{a}} - \alpha \mathbf{p}_1|^2}{2\beta^2}\right) \\ &= \sqrt{\frac{1}{(2\pi)^n \beta^{2l}}} \exp\left(-\frac{\mathbf{h}^T (\mathbb{1} - Q_l Q_l^T) \mathbf{h}}{2} - \frac{(\mathbf{h} - \alpha \mathbf{q}_1)^T Q_l Q_l^T (\mathbf{h} - \alpha \mathbf{q}_1)}{2\beta^2}\right). \end{aligned} \quad (\text{A.4})$$

It is worth noting that $\mathcal{L}(\mathbf{h}|\alpha, \beta, l)$ is a multivariate normal distribution having

$$\langle \mathbf{h} \rangle_{\mathbf{h}|\alpha, \beta, l} = \alpha \mathbf{q}_1 \quad (\text{A.5})$$

mean and

$$\Sigma_h = \mathbb{1} + (\beta^2 - 1) Q_l Q_l^T \quad (\text{A.6})$$

variance-covariance matrix. To prove (A.5), it is enough to find the $\mathcal{L}(\mathbf{h}|\alpha, \beta, l)$ maximum by observing that the gradient of the quadratic form $\mathbf{h}^T \Sigma \mathbf{h}$, where Σ is symmetric, is $\partial_{\mathbf{h}}(\mathbf{h}^T \Sigma \mathbf{h}) = 2\mathbf{h}^T \Sigma$ and that $\mathbf{q}_1^T Q_l Q_l^T = \mathbf{q}_1^T$ [14].

Next, we observe that $\partial_{\mathbf{h}}^2(\mathbf{h}^T \Sigma \mathbf{h}) = 2\Sigma$. Hence, the opposite of the Hessian of the $\mathcal{L}(\mathbf{h}|\alpha, \beta, l)$ exponent,

$$\Sigma_h^{-1} = \mathbb{1} - \frac{\beta^2 - 1}{\beta^2} Q_l Q_l^T, \quad (\text{A.7})$$

is the inverse of (A.6) [24, 25].

Appendix B. Fisher information matrix

The entries of the Fisher information matrix (14) are

$$F_{\alpha\alpha} = -\langle \partial_{\alpha\alpha} [\ln(\mathcal{L})] \rangle_{h|\alpha,\beta,l} = |\mathbf{p}_1|^2 / \beta^2 = 1 / \beta^2, \quad (\text{B.1})$$

$$F_{\alpha\beta} = -\langle \partial_{\alpha\beta} [\ln(\mathcal{L})] \rangle_{h|\alpha,\beta,l} = \langle \hat{\mathbf{a}} - \alpha \mathbf{p}_1 \rangle_{h|\alpha,\beta,l} \mathbf{p}_1 / \beta^3 = 0, \quad (\text{B.2})$$

$$F_{\beta\beta} = -\langle \partial_{\beta\beta} [\ln(\mathcal{L})] \rangle_{h|\alpha,\beta,l} = \frac{1}{\beta^2} \left(\frac{3}{\beta^2} \langle |\hat{\mathbf{a}} - \alpha \mathbf{p}_1|^2 \rangle_{h|\alpha,\beta,l} - l \right) = \frac{2l}{\beta^2}, \quad (\text{B.3})$$

where, apart from unessential terms independent of α and β ,

$$\ln(\mathcal{L}) = -|\hat{\mathbf{a}} - \alpha \mathbf{p}_1|^2 / (2\beta^2) - l \ln(\beta). \quad (\text{B.4})$$

To evaluate $\langle \hat{\mathbf{a}} - \alpha \mathbf{p}_1 \rangle_{h|\alpha,\beta,l}$, we observe that $\hat{\mathbf{a}} = Q_l^T \mathbf{h}$ and use (A.5) and $Q_l^T \mathbf{q}_1 = \mathbf{p}_1$. The evaluation of $\langle |\hat{\mathbf{a}} - \alpha \mathbf{p}_1|^2 \rangle_{h|\alpha,\beta,l}$ is a bit more tricky. First, we note that $\hat{\mathbf{a}} = Q_l^T \mathbf{h}$ and $\langle \hat{\mathbf{a}} \rangle_{h|\alpha,\beta,l} = Q_l^T \langle \mathbf{h} \rangle_{h|\alpha,\beta,l} = \alpha Q_l^T \mathbf{q}_1 = \alpha \mathbf{p}_1$. Hence,

$$\langle |\hat{\mathbf{a}} - \alpha \mathbf{p}_1|^2 \rangle_{h|\alpha,\beta,l} = \text{Tr}(\Sigma_{\hat{\mathbf{a}}}), \quad (\text{B.5})$$

where $\Sigma_{\hat{\mathbf{a}}}$ is the variance-covariance matrix of $\hat{\mathbf{a}}$. Next, since $\hat{\mathbf{a}} = Q_l^T \mathbf{h}$, by using $Q_l^T Q_l = \mathbb{1}$ and (A.6),

$$\Sigma_{\hat{\mathbf{a}}} = Q_l^T \Sigma_h Q_l = \beta^2 \mathbb{1}. \quad (\text{B.6})$$

Eventually, $\text{Tr}(\Sigma_{\hat{\mathbf{a}}}) = l\beta^2$.

Appendix C. Marginal likelihood

By carrying out the needed integrations with the aid of Mathematica[®] [17], The marginal likelihood is

$$\begin{aligned} Z(\mathbf{h}|l) &= \int_1^{+\infty} \int_{-\infty}^{+\infty} \mathcal{L}(\mathbf{h}|\alpha, \beta, l) \pi_J(\alpha, \beta) d\alpha d\beta \\ &\propto \exp\left(-\frac{|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2}{2}\right) \int_1^{+\infty} \left[\frac{1}{\beta^{l+2}} \exp\left(\frac{|\hat{\mathbf{a}}'|^2}{2\beta^2}\right) \int_{-\infty}^{+\infty} \exp\left(\frac{(\hat{a}_1 - \alpha)^2}{2\beta^2}\right) d\alpha \right] d\beta \\ &\propto \exp\left(-\frac{|\hat{\mathbf{e}}|^2}{2}\right) \int_1^{+\infty} \frac{1}{\beta^{l+1}} \exp\left(\frac{|\hat{\mathbf{a}}'|^2}{2\beta^2}\right) d\beta \\ &= \frac{\sqrt{2^{l-1}} \exp(-|\hat{\mathbf{e}}|^2/2) \gamma(l/2, |\hat{\mathbf{a}}'|^2/2)}{|\hat{\mathbf{a}}'|^l}, \end{aligned} \quad (\text{C.1})$$

where $\gamma(s, z_2)$ is the lower incomplete gamma function [16], $|\hat{\mathbf{e}}|^2 = \hat{\mathbf{e}}^T \hat{\mathbf{e}}$ is the sum of the squared residuals, $\hat{\mathbf{e}} = \mathbf{h} - \hat{\mathbf{h}}$ are the residuals, and $\hat{\mathbf{h}} = Q_l \hat{\mathbf{a}}$ is the least-squares estimate of \mathbf{h} . To prove that $|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2 = |\hat{\mathbf{e}}|^2$, which has been used to obtain (C.1), we observe that $\hat{\mathbf{h}}^T (\mathbf{h} - \hat{\mathbf{h}}) = 0$, because $\hat{\mathbf{h}}$ is the projection of \mathbf{h} in the $Q_l \mathbf{a}$ sub-space. Furthermore, $|\hat{\mathbf{a}}|^2 = \mathbf{h}^T Q_l Q_l^T \mathbf{h} = \mathbf{h}^T \hat{\mathbf{h}}$. Hence,

$$|\mathbf{h}|^2 - |\hat{\mathbf{a}}|^2 = \mathbf{h}^T (\mathbf{h} - \hat{\mathbf{h}}) = (\mathbf{h} - \hat{\mathbf{h}})^T (\mathbf{h} - \hat{\mathbf{h}}) = \hat{\mathbf{e}}^T \hat{\mathbf{e}}. \quad (\text{C.2})$$

Since $Q^T \mathbf{h}' = Q^T \mathbf{h}$ and $Q Q^T \mathbf{q}_1 = \mathbf{q}_1$, the marginal likelihood (C.1) is invariant against the $\mathbf{h}' = \mathbf{h} + g \mathbf{q}_1$ transformation, where g is any numerical constant. This can be used to rewrite it in a form that is independent of the data normalization and the orthonormalization of the design matrix.

To this end, firstly, we shift the \mathbf{y} data in such a way that their generalised mean (9a-b) is null. In such a way $\hat{a}_1 = Q_1^T \mathbf{h}$ shifts to zero. Since $|\mathbf{q}_1| = Q_1^T \mathbf{q}_1 = 1$,

this corresponds to shift \mathbf{h} to $\mathbf{h} - (\mathbf{q}_1^T \mathbf{h}) \mathbf{q}_1$ and keeps the $Z(\mathbf{h}|l)$ value unchanged. Now, when $\hat{a}_1 = 0$, we have $|\hat{\mathbf{a}}'|^2 = |\hat{\mathbf{a}}|^2$. In addition, by using $\mathbf{Q}_l^T \mathbf{Q}_l = \mathbb{1}$, we get $|\hat{\mathbf{a}}|^2 = \hat{\mathbf{a}}^T (\mathbf{Q}_l^T \mathbf{Q}_l) \hat{\mathbf{a}} = |\hat{\mathbf{h}}|^2$. Eventually, in (C.1), we have $|\hat{\mathbf{a}}'|^2 = |\hat{\mathbf{h}}|^2$.

Secondly, we observe that $|\hat{\mathbf{h}}|^2 = \hat{\mathbf{y}}^T \Sigma_u^{-1} \hat{\mathbf{y}} = |\hat{\mathbf{y}}|^2$ and $|\hat{\epsilon}|^2 = \hat{\mathbf{u}}^T \Sigma_u^{-1} \hat{\mathbf{u}} = |\hat{\mathbf{u}}|^2$, where $\hat{\mathbf{y}}$ are the generalised least-squares estimate of the data (17) after subtracting their generalised mean and $\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}}$.

To conclude, we can rewrite (C.1) by substituting $|\mathbf{u}|$ and $|\hat{\mathbf{y}}|$ for $|\hat{\epsilon}|$ and $|\hat{\mathbf{a}}'|$, so that (16) gives the marginal likelihood in terms of the original data and variance-covariance matrix.

Appendix D. Design matrices

The central second moment of a paraxial beam depends quadratically on the propagation distance, z . Therefore, after assembling the data – three series of seven measured values – in the vector $\mathbf{y} = (y_{11}, \dots, y_{17}, y_{21}, \dots, y_{27}, y_{31}, \dots, y_{37})^T$ and the nine parameters of three non-intersecting parabolas in $\mathbf{b} = (a_1, b_1, c_1, a_2, b_2, c_2, a_3, b_3, c_3)^T$, the design matrix is

$$W_9 = \begin{pmatrix} 1 & z_1 & z_1^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ \dots & & & & & & & & \\ 1 & z_7 & z_7^2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & z_1 & z_1^2 & 0 & 0 & 0 \\ \dots & & & & & & & & \\ 0 & 0 & 0 & 1 & z_7 & z_7^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & z_1 & z_1^2 \\ \dots & & & & & & & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & z_7 & z_7^2 \end{pmatrix}. \quad (\text{D.1})$$

A quartic polynomial was used as an alternative explanation to test the measurement quality. The design matrix,

$$W_{12} = \begin{pmatrix} 1 & z_1 & z_1^2 & z_1^4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \dots & & & & & & & & & & & \\ 1 & z_7 & z_7^2 & z_7^4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & z_1 & z_1^2 & z_1^4 & 0 & 0 & 0 & 0 \\ \dots & & & & & & & & & & & \\ 0 & 0 & 0 & 0 & 1 & z_7 & z_7^2 & z_7^4 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & z_1 & z_1^2 & z_1^4 \\ \dots & & & & & & & & & & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & z_7 & z_7^2 & z_7^4 \end{pmatrix}, \quad (\text{D.2})$$

is obtained from (D.1) by adding three columns related to the z^4 basis and by supplementing the parameter vector, $\mathbf{b} = (a_1, b_1, c_1, d_1, a_2, b_2, c_2, d_2, a_3, b_3, c_3, d_3)^T$, accordingly.

Eventually, the expectation that the three parabolas intersect in the null-defocus configuration was modelled by imposing the constraint $a_1 = a_2 = a_3 = a_0$. This was

done by using the design matrix

$$W_7 = \begin{pmatrix} 1 & z_1 & z_1^2 & 0 & 0 & 0 & 0 \\ \dots & & & & & & \\ 1 & z_7 & z_7^2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & z_1 & z_1^2 & 0 & 0 \\ \dots & & & & & & \\ 1 & 0 & 0 & z_7 & z_7^2 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & z_1 & z_1^2 \\ \dots & & & & & & \\ 1 & 0 & 0 & 0 & 0 & z_7 & z_7^2 \end{pmatrix}, \quad (\text{D.3})$$

where the seven parameters are $\mathbf{b} = (a_0, b_1, c_1, b_2, c_2, b_3, c_3)^T$.

References

- [1] Birge R T 1932 *Phys. Rev.* **40**(2) 207–227
- [2] Mana G, Massa E and Predescu M 2012 *Metrologia* **49** 492–500
- [3] Wuebbeler G, Bodnar O and Elster C 2016 *Metrologia* **53** 1131–1138
- [4] Gull G F 1988 Bayesian inductive inferences and maximum entropy *Maximum entropy and Bayesian meyhods in science and engineering* ed Erickson G and Smith C (Dordrecht: Kluwer Academic Publishers) pp 53–74
- [5] Nesseris S and García-Bellido J 2013 *Journal of Cosmology and Astroparticle Physics* **2013** 036
- [6] Mana G, Albo P A G and Lago S 2014 *Measurement* **55** 564–570
- [7] Amari S, Nagaoka H and Harada D 2007 *Methods of Information Geometry* Translations of Mathematical Monographs vol. 191 (Oxford: Oxford University Press)
- [8] Arwini K and Dodson C 2008 *Information Geometry: Near Randomness and Near Independence* Lecture Notes in Mathematics 1953 (Berlin Heidelberg: Springer)
- [9] Jeffreys H 1946 *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* **186** 453–461
- [10] Jeffreys H 1998 *The Theory of Probability* (Oxford: Oxford University Press)
- [11] Kass R E and Wasserman L 1996 *Journal of the American Statistical Association* **91** 1343–1370
- [12] Costa S I, Santos S A and Strapasson J E 2014 *Discrete Applied Mathematics* **197** 59–69
- [13] Golub G and Van Loan C 1996 *Matrix Computations* Johns Hopkins Studies in the Mathematical Sciences (Johns Hopkins University Press)
- [14] Dhrymes P 1978 *Mathematics for Econometrics* (Berlin Heidelberg: Springer-Verlag)
- [15] Berger J O, Bernardo J M and Sun D 2009 *Ann. Statist.* **37** 905–938
- [16] Weisstein E W Incomplete gamma function
<http://mathworld.wolfram.com/IncompleteGammaFunction.html>
from MathWorld – A Wolfram Web Resource Accessed: 2017-01-12
- [17] Wolfram Research Inc 2017 Mathematica 11.2
- [18] Akaike H 1974 *IEEE Transactions on Automatic Control* **19** 716–723
- [19] Burnham K and Anderson D 2002 *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach* (New York: Springer)
- [20] Schwarz G 1978 *Ann. Statist.* **6** 461–464
- [21] Konishi S and Kitagawa G 2008 *Information Criteria and Statistical Modeling* Springer series in statistics (New York: Springer)
- [22] Bhat H and Kumar N 2010 On the derivation of the bayesian information criterion Tech. rep. School of Natural Sciences, University of California, Merced
- [23] Mana G, Massa E, Sasso C P, Andreas B and U K 2017 *Metrologia* **54** 559–565
- [24] Henderson H V and Searle S R 1981 *SIAM Review* **23** 53–60
- [25] Tylavsky D J and Sohie G R L 1988 *Proceedings of the IEEE* **74** 1050–52