Bayesian Unbiasing of the Gaia space mission time series database

Héctor E. Delgado and Luis M. Sarro

Dpto. de Inteligencia Artificial , UNED, Juan del Rosal, 16, 28040 Madrid, Spain hed_up@iasystems.org lsb@dia.uned.es

Abstract. 21^{st} century astrophysicists are confronted with the herculean task of distilling the maximum scientific return from extremely expensive and complex space- or ground-based instrumental projects. This paper concentrates in the mining of the time series catalog produced by the European Space Agency Gaia mission, launched in December 2013. We tackle in particular the problem of inferring the true distribution of the variability properties of Cepheid stars in the Milky Way satellite galaxy known as the Large Magellanic Cloud (LMC). Classical Cepheid stars are the first step in the so-called distance ladder: a series of techniques to measure cosmological distances and decipher the structure and evolution of our Universe. In this work we attempt to unbias the catalog by modelling the aliasing phenomenon that distorts the true distribution of periods. We have represented the problem by a 2-level generative Bayesian graphical model and used a Markov chain Monte Carlo (MCMC) algorithm for inference (classification and regression). Our results with synthetic data show that the system successfully removes systematic biases and is able to infer the true hyperparameters of the frequency and magnitude distributions.

Keywords: Astrostatistics, Bayesian, data analysis, hierarchical model, Markov chain Monte Carlo, catalogues

1 Introduction

Gaia [5] is a European Space Agency (ESA) space mission, launched in December 2013, whose main objective is to compile a large-scale astronomical survey of about one billion stars ($\approx 1\%$) of our Galaxy and its Local Group. The satellite will scan the entire sky for about 5 years yielding an unprecedented catalog in both size and precision of positions, distances and proper motion measures. Additionally, it will perform multi-epoch photometry (70 transits per object on average) which renders the satellite suitable too for studies of stellar variability. Amongst the many variability types present in the stellar zoo, one in particular is of paramount importance: the Classical Cepheids. Classical Cepheids represent the first calibrator in the cosmic distance ladder used to infer the structure and evolution of our Universe, and our current knowledge about the Big Bang,

the inflationary period, the dark matter problem or dark energy relies on the period-luminosity relation for Classical Cepheids [14]. Therefore, a precise and accurate understanding of the population of Classical Cepheids is central to all cosmological studies.

In this paper we address the problem of inferring the true properties of this population of variable stars from the petabyte-size *Gaia* catalog. In order to populate the catalog, and as part of a much larger framework to deliver a data set of scientific quality, the Data Processing and Analysis Consortium (DPAC¹) developed a pipeline to characterize the time series observed and classify them. A key element of this process is that the time sampling of stellar brightness time series will have the imprint of the satellite intrinsic frequencies (amongst other, the spinning and precessing frequencies, a description of which is out of the scope of this paper). As a consequence, some (but not all) of the derived frequencies will be affected by aliasing which results in biased samples.

The objective is to characterize the phenomenon of *aliasing* in the *Gaia* catalog, correct for it, and reconstruct the real distribution of LMC Classical Cepheids properties. In order to achieve these goals, we tackle the problem under the Bayesian paradigm [6,7] and adopt the knowledge representation language of Bayesian Networks (BN) [11,8]. This framework allows a hierarchical representation of the problem in which the time series gathered by Gaia are the product of a generative process which ultimately depends on the parameters of the population of stars. Given that the computation of the posterior probabilities of our model are analytically intractable, the inference mechanism of our proposal is founded in Markov chain Monte Carlo (MCMC) simulation techniques [13].

We have validated our models using a data base of 36688 synthetic Classical LMC Cepheids time series generated according to controlled prescriptions based on current understanding of the true distributions and the satellite characteristics. Our results prove that we are ready for the second *Gaia* data release expected for 2018. This will be the first data release to include photometric time series (although this still needs to be confirmed).

The structure of the rest of the paper is as follows. In section 2 we describe our model and the MCMC technique used for the inference of the parameters of interest. In Section 3 we validate the model with the simulated data base in a scenario of extreme aliasing and describe the results of this validation procedure. Finally, in Section 4 we summarize the contributions of this work and some of its limitations, and give pointers to future developments.

2 Hierarchical Modelling of the distribution of pulsation properties of Classical Cepheid Variable stars

2.1 The Hierarchical Model

Figure 1 and Table 1 depict the structure of the DAG associated to the model and summarize the meanings of the nodes and the types of their distributions.

¹ The DPAC (Data Processing and Analysis Consortium) is the consortium responsible for building and making accessible the GAIA catalogue.

We classify the nodes into a hierarchy of three levels. The hierarchy distinguishes between evidential nodes (observations), the rest of nodes inside the rectangle or *plate*, which is replicated N times (one per star), and the nodes outside the rectangle. In the following paragraphs we describe the parameters and probability distributions for each level and its contribution to the joint probability distribution.



Fig. 1. Graph structure of our proposed Bayesian Graphical Model (BGM). Most fixed parameters are not included in the graph, with the exception of those enclosed inside a square. See the text and Table 1 for node descriptions.

2.1.1 Likelihood. In the bottom level of our graph we present the *evidential nodes*, that is, the variables measured directly or derived by the DPAC

$$\mathcal{D} = (\nu_{\text{rec},i}, A_{\text{rec},i}, m_{G_{\text{rec}},i}) \ . \tag{1}$$

These nodes, depicted by double circles, are the output/recovered frequency $\nu_{\text{rec},i}$, the amplitude $A_{\text{rec},i}$ and the apparent G-magnitude $m_{G_{\text{rec}},i}$ for the *i*-th star.

Recovered Frequencies. Most of the pairs (ν_{input}, ν_{rec}) in the simulated data base fall on straights lines of the form:

$$\nu_{\rm rec} = \pm \nu_{\rm input} \pm k_1 \nu_s \pm k_2 \nu_p , \qquad (2)$$

where $k_1 \in \{0, 3, 7\}, k_2 \in \{0, ..., 19\}, \nu_s \approx \frac{1}{0.25} = 4d^{-1}$ is the rotational frequency of Gaia and $\nu_p = \frac{1}{63}d^{-1}$ is its precessional frequency. We refer to each line as a

4

Node	Description	Type of distribution	
$ au_G$	Precision	Gamma NI prior	
\mathbf{a}_G	Slopes	Gaussian NI prior	
$oldsymbol{b}_G$	Intercepts	Gaussian NI prior	
$m_{G,i}$	Input apparent G magnitude	Gaussian	
$m_{G,\mathrm{rec},i}$	Recovered apparent G magnitude	Gaussian	
μ_A	Mean	Gaussian NI prior	
$ au_A$	Precision	Gamma NI prior	
a_A	Slope	Gaussian NI prior	
b_A	Intercept	Gaussian NI prior	
A_i	Input amplitude	Gaussian	
$A_{\mathrm{rec},i}$	Recovered amplitude	Mixture of skewed Cauchy	
$w_ u$	Mixing proportions.	Gamma NI prior	
T_{ν_i}	Category of $\log(\nu_i)$	Categorical	
$\mu_{ u}$	Mean	Non informative	
$oldsymbol{ heta}_{ u}$	Mean Perturbations	Gaussian NI prior	
$ au_ u$	Precision	Non informative	
$oldsymbol{\omega}_ u$	Precision Perturbations	Uniform prior	
$\log(\nu_i)$	Input frequency $\left[d^{-1}\right]$.	Mixture of Gaussian	
Λ	Logistic R. coefficients	Student t prior	
$T_{\nu_{\mathrm{rec},i}}$	Category of $\nu_{\mathrm{rec},i}$	Categorical	
$\nu_{\mathrm{rec},i}$	Recovered frequency	Mixture of Gaussian	

Table 1. Description of parameters. NI = non informative

locus/*category* of recovered frequencies. Excluding the line $\nu_{\rm rec} = \nu_{\rm input}$, all these *loci* correspond to spurious (aliased) frequencies. Based on that, we parameterize the *i*-th recovered frequency as the following mixture of Gaussian distributions

$$f\left(\nu_{\text{rec},i} \mid \log\left(\nu_{i}\right), T_{\nu_{\text{rec},i}}\right) = \sum_{j=1}^{M} \delta_{T_{\nu_{\text{rec},i}}}^{j} \mathsf{N}\left((-1)^{j-1} \, 10^{\log(\nu_{i})} + b_{j}, \tau_{\nu_{\text{rec}}}\right) \,.$$
(3)

In Equation 3 the Kronecker deltas $\delta^j_{T_{\nu_{\text{rec},i}}}$ dictate the Gaussian component to which $\nu_{\text{rec},i}$ belongs according to the value of the categorical variable $T_{\nu_{\text{rec},i}}$ (described in Section 2.1.2). The mean of each component represents the *locus* in which the input frequency has been recovered, i.e. the identity locus, with $b_j = 0$ for j = 1, or some locus of spurious (aliased) frequencies for j > 1. We assume the same precision $\tau_{\nu_{\text{rec}}} = 10000$ for all components.

Recovered Amplitudes. To gain insight into the form of the conditional distribution of the recovered amplitude given the input amplitude we have checked the hypothesis that recovered amplitudes are also biased by the aliasing phenomenon, just as recovered frequencies are. By analysing the relationship between *loci* of frequencies and pairs (A_{input}, A_{rec}), we have discovered that for a perfect recovery the distribution $A_{rec} \mid A$ is skewed to lower amplitudes with a central parameter approximately equal to the input amplitude. Otherwise, for *loci* of aliased frequencies we have observed that the skewness of the recovered

amplitude increases as the input amplitude does according to a certain slope to be determined as part of the model. To account for this fact we have fitted two linear regression models

$$A_{\text{rec},i} = \beta_1^j A_{\text{in},i} + \beta_0^j + \epsilon_i^j, \ j = 1,2 \ , \tag{4}$$

with j = 1 corresponding to the identity locus and j = 2 to the *loci* $\nu_{\rm rec} = \pm \nu_{\rm in} + 7\nu_s - 3\nu_p^{-2}$. For the identity locus, we have assumed a skewed Student t distribution [2] with one degree of freedom (skewed Cauchy) for the error component $\epsilon_i^1 \sim \operatorname{st}(0, \omega, \alpha, 1)$ where ω and α denote respectively the shape and scale parameters. For the *locus* $\nu_{\rm rec} = \pm \nu_{\rm in} + 7\nu_s - 3\nu_p$ we have assumed that $\epsilon_i^2 \sim \operatorname{t}(0, \omega, 1)$. Based on that, we model the conditional distribution for the recovered amplitude $A_{{\rm rec},i}$ by means of the mixture of two skewed Student t distributions

$$f\left(A_{\text{rec},i} \mid A_i, T_{\nu_{\text{rec},i}}\right) = \delta^{1}_{T_{\nu_{\text{rec},i}}} \mathsf{ST}\left(A_i, 0.020, -2.395, 1\right) \\ + \sum_{j=2}^{M} \delta^{j}_{T_{\nu_{\text{rec},i}}} \mathsf{ST}\left(0.749 \cdot A_i, 0.0266, 0, 1\right) ,$$
(5)

where the location parameters $\xi_1 = A_i$, $\xi_j = 0.749 \cdot A_i$, $\forall j = 2, ..., M$, the scale ω and the shapes α have been obtained from the fitting of the two linear models of Equation 4 and taken as constants in our BGM.

Recovered Apparent Magnitudes. We parameterize the distribution of the *i*-th recovered apparent G magnitude by means of a Gaussian distribution with mean $m_{G,i}$ and precision $\tau_{G(rec)} = 2.5\text{E}+5$ (to be adjusted when real Gaia data become available)

$$f\left(m_{G_{rec},i} \mid m_{G,i}\right) = \mathsf{N}\left(m_{G,i}, \tau_{G_{rec}}\right) \ . \tag{6}$$

The conditional distribution of the data given their parents is then given by

$$p\left(\mathcal{D} \mid \boldsymbol{\theta}_{1}\right) = \prod_{i=1}^{N} f_{1}\left(\nu_{\mathrm{rec},i} \mid \log\left(\nu_{i}\right), T_{\nu_{\mathrm{rec},i}}\right) \cdot f_{2}\left(A_{\mathrm{rec},i} \mid A_{i}, T_{\nu_{\mathrm{rec},i}}\right)$$
(7)

$$\cdot f_{3}\left(m_{G_{\mathrm{rec},i}} \mid m_{G,i}\right) .$$

2.1.2 First Level Random Parameters. These are

$$\boldsymbol{\theta}_1 = \left(\log\left(\nu_i\right), A_i, m_{G,i}, T_{\nu_{\text{rec},i}}, T_{\nu_i} \right) \ . \tag{8}$$

In θ_1 , we distinguish two classes of nodes. The *input nodes* are, for the *i*-th star, the real frequency $\log(\nu_i)$, the real amplitude A_i and the real apparent G-magnitude $m_{G,i}$. The *categorical nodes* $T_{\nu_{rec,i}}$ and T_{ν_i} determine the component

 $^{^{2}}$ We only select these particular *loci* of aliased frequencies because they are the most frequent *loci* located far away from the identity *locus* and because the model does not work well if we include more *loci* located close to them.

of a node modelled by a mixture of distributions. T_{ν_i} and $T_{\nu_{\text{rec},i}}$ are respectively associated with the real frequency and the recovered frequency and amplitude. In Figure 1 all the nodes at this level replicate with the plate. They depend on (amongst other) non informative orphan nodes outside the plate.

Categories of Recovered Frequencies. The node $T_{\nu_{\text{rec},i}}$ takes a value $j \in \{1, ..., M\}$ if the *i*-th frequency has been recovered in the *j*-th locus, which occurs with a probability π_{ij} . In this paper we assume that the main factor determining the aliasing phenomenon in Gaia is the ecliptic latitude β of the stars. The influence of β over the rate of correct detections of periodic signals by *Gaia* has been studied in [4] where it is shown that for high values of β , typical of LMC sources, the relation between the rate of correct detections and β is approximately linear with a negative slope. Based on that, we make π_{ij} depend on the ecliptic latitude β_i and parameterize this dependence by a multinomial logistic regression submodel with a *softmax* transfer function. We model the conditional distribution of $T_{\nu_{\text{rec},i}}$ as

$$p\left(T_{\nu_{\mathrm{rec},i}} \mid \{\boldsymbol{\lambda}_j\}_{j=2}^M\right) = \mathsf{Cat}\left(M, \{\pi_{ij}\left(\beta'_i, \boldsymbol{\lambda}_j\right)\}_{j=1}^M\right) , \qquad (9)$$

with

$$\pi_{ij}\left(\beta_i', \boldsymbol{\lambda}_j\right) = \frac{e^{\boldsymbol{\lambda}_j^T \cdot \left(1, \beta_i'\right)}}{\sum_{l=1}^M e^{\boldsymbol{\lambda}_l^T \cdot \left(1, \beta_i'\right)}} , \qquad (10)$$

where we have rescaled the predictor β_i by subtracting the mean and dividing by two times the standard deviation, i.e. $\beta'_i = \frac{\beta_i - \overline{\beta}}{2 \cdot \operatorname{sd}(\beta)}$, which guaranties that the mean and the standard deviation are respectively 0 and 0.5.

Input Frequencies and Categories. The marginal distribution of the (decadic) logarithm of the input frequency in the synthetic data set created by the DPAC Quality Assessment group was sampled from a mixture of five Gaussian distributions [1]. In our BGM, we parameterize it by the mixture of only three components³

$$f\left(\log\left(\nu_{i}\right) \mid T_{\nu_{i}}, \mu_{\nu}, \boldsymbol{\theta}_{\nu}, \tau_{\nu}, \boldsymbol{\omega}_{\nu}\right) = \delta_{T_{\nu_{i}}}^{1} \mathsf{N}\left(\mu_{\nu}, \tau_{\nu}\right) + \delta_{T_{\nu_{i}}}^{2} \mathsf{N}\left(\mu_{\nu} + \sqrt{\tau_{\nu}^{-1}} \theta_{\nu 1}, \tau_{\nu} \omega_{\nu 1}^{-2}\right) + \delta_{T_{\nu_{i}}}^{3} \mathsf{N}\left(\mu_{\nu} + \sqrt{\tau_{\nu}^{-1}} \theta_{\nu 1} + \sqrt{\tau_{\nu}^{-1}} \omega_{\nu 1} \theta_{\nu 2}, \tau_{\nu} \omega_{\nu 1}^{-2} \omega_{\nu 2}^{-2}\right)$$
(11)

In Equation 11, μ_{ν} and τ_{ν} denote, respectively, the mean and the precision of the first component of the mixture. $(\theta_{\nu 1}, \theta_{\nu 2})$ and $(\omega_{\nu 1}, \omega_{\nu 2})$ denote, respectively, the perturbation parameters which affect the mean and the scale parameter of a given component to obtain the mean and scale parameter of the next component [12]. The Kronecker deltas $\delta_{T_{\nu_i}}^j$ have the same role as in Eq. 3 but now the

 $^{^{3}}$ We rely on the Occam's razor principle.

categorical variable T_{ν_i} represents the class of the real frequency. For T_{ν_i} we assign the distribution

$$p(T_{\nu_i}) = \mathsf{Cat}(3, w_{\nu 1}, w_{\nu 2}, w_{\nu 3}) , \qquad (12)$$

where $w_{\nu j}$ are the mixing proportions of the mixture.

Input Amplitudes. This distribution has been simulated based on the OGLE III catalogue of Classical Cepheids [15], as

$$f(A \mid \log(\nu)) = \begin{cases} \mathsf{N}(-0.5 \cdot \log(\nu) + 0.2, 0.15) & \log(\nu) < -1\\ \mathsf{N}(0.7, 0.15) & \log(\nu) > -1 \end{cases}$$
(13)

In our BGM we parameterize this variable as

$$f(A_i \mid \log(\nu_i), a_A, b_A, \mu_A, \tau_A) = \mathbf{1}_{\{\log(\nu_i) < -1\}} \mathsf{N}(a_A \cdot \log(\nu_i) + b_A, \tau_A) + \mathbf{1}_{\{\log(\nu_i) > -1\}} \mathsf{N}(\mu_A, \tau_A) ,$$
(14)

where $\mathbf{1}_S$ denotes the indicator function of a subset S, a_A and b_A are, respectively, the slope and the intercept of the regression line of A on $\log(\nu)$ when $\log(\nu) < -1$, μ_A denotes the mean of the amplitude when $\log(\nu) > -1$, and τ_A denotes the precision, which we take equal in both cases.

Input Apparent G magnitudes. Based on Equations 12 and 13 of [14] and discarding the distance r to the sources, we parameterize this node as

$$f(m_{G,i} | \log(\nu_i), a_{G1}, b_{G1}, a_{G2}, b_{G2}, \tau_G) = \mathbf{1}_{\{\log(\nu_i) < -1\}} \mathsf{N}(a_{G1} \cdot \log(\nu_i) + b_{G1}, \tau_G)$$
(15)
+ $\mathbf{1}_{\{\log(\nu_i) > -1\}} \mathsf{N}(a_{G2} \cdot \log(\nu_i) + b_{G2}, \tau_G)$.

The conditional distribution of the first level of random parameters given the parameters of the top level is then

$$p(\boldsymbol{\theta}_{1} \mid \boldsymbol{\theta}_{2}) = \prod_{i=1}^{N} g_{1} \left(T_{\nu_{rec,i}} \mid \{\boldsymbol{\lambda}_{j}\}_{j=2}^{M} \right) \cdot g_{2} \left(A_{i} \mid \log\left(\nu_{i}\right), a_{A}, b_{A}, \mu_{A}, \tau_{A} \right)$$

$$\cdot g_{3} \left(m_{G,i} \mid \log\left(\nu_{i}\right), \mathbf{a}_{G}, \mathbf{b}_{G}, \tau_{G} \right) \cdot g_{4} \left(\log\left(\nu_{i}\right) \mid T_{\nu_{i}}, \lambda_{\nu}, \boldsymbol{\theta}_{\nu}, \tau_{\nu}, \boldsymbol{\omega}_{\upsilon} \right)$$

$$\cdot g_{5} \left(T_{\nu_{i}} \mid \boldsymbol{w}_{\nu} \right)$$

$$(16)$$

2.1.3 Top Level Random Parameters. These hyperparameters are

$$\boldsymbol{\theta}_2 = (a_A, b_A, \mu_A, \tau_A, \mathbf{a}_G, \mathbf{b}_G, \tau_G, \mu_\nu, \boldsymbol{\theta}_\nu, \tau_\nu, \boldsymbol{\omega}_\nu, \boldsymbol{w}_\nu, \Lambda) \quad . \tag{17}$$

 θ_2 include the orphan nodes in the graph. We only have a vague (or non informative) prior knowledge about their distributions. The nodes denoted by a and b represent the slopes and intercepts of the distributions of the real amplitude and apparent G-magnitude given the frequency. The nodes denoted by

 τ and μ represent precisions and means. The nodes denoted by Λ represent the coefficients of the logistic regression submodel of Equation 10. The rest of nodes are associated with the parameterization of the real frequency of Equation 11. For these latter hyperparameters we take the non informative priors

$$p(\boldsymbol{w}_{\nu}) = \mathsf{Dir}(1, 1, 1) \tag{18}$$

$$p(\mu_{\nu}) = \mathsf{N}(0, 0.001) \tag{19}$$

$$p\left(\theta_{\nu j}\right) = \mathsf{N}\left(0, 0.01\right) \tag{20}$$

$$p(\tau_{\nu}) = \text{Gamma}(0.001, 0.001) \tag{21}$$

$$p\left(\omega_{\nu j}\right) = \mathsf{U}\left(0,1\right) \tag{22}$$

For the hyperparameters of the logistic regression submodel of Equation 10 $\lambda_j = (\lambda_{0j}, \lambda_{1j})$ with $j \in \{2, ..., M\}$, we assign the weakly informative priors $p(\lambda_{kj}) = t(0, \frac{1}{2.5^2}, 7), k \in \{0, 1\}$. This election provides a minimal prior information to constrain the range of coefficients λ_{kj} once the covariate β_i has been rescaled [6]. This approximation is used to enhance the convergence rate of our model.

For the parameters a_A , b_A , λ_A of the input amplitude distribution of Equation 14 and the parameters a_{G1} , b_{G1} , a_{G2} , b_{G2} of the input apparent G magnitude of Equation 15 we take N (0,0.001) non informative priors. And for the precisions τ_A and τ_G we take Gamma (0.001, 0.001) priors. For all these priors the full conditional distribution of the node is available in closed form.

The distribution (hyperprior) of the top level parameters is then

$$p(\boldsymbol{\theta}_2) = h_1(a_A) \cdot h_2(b_A) \cdot h_3(\mu_A) \cdot h_4(\tau_A) \cdot h_5(\mathbf{a}_G) \cdot h_6(\mathbf{b}_G) \cdot h_7(\tau_G)$$

$$\cdot h_8(\boldsymbol{w}_{\nu}) \cdot h_9(\mu_{\nu}) \cdot h_{10}(\boldsymbol{\theta}_{\nu}) \cdot h_{11}(\tau_{\nu}) \cdot h_{12}(\boldsymbol{\omega}_{\nu}) \cdot h_{13}(\Lambda) .$$
(23)

2.1.4 Joint distribution of the Parameters and Data. From Equations 7, 16 and 23 we formulate the joint PDF associated to the graphical mode by

$$p(\boldsymbol{\theta}, \mathcal{D}) = p(\mathcal{D} \mid \boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}) = p(\mathcal{D} \mid \boldsymbol{\theta}_1) \cdot p(\boldsymbol{\theta}_1 \mid \boldsymbol{\theta}_2) \cdot p(\boldsymbol{\theta}_2) .$$
(24)

2.2 Computation

The joint posterior distribution of the 22+5N parameters of the model described in Section 2.1 is given by

$$\pi^* (\boldsymbol{\theta}) = \pi (\boldsymbol{\theta} \mid \mathcal{D}) \propto \mathcal{L} (\boldsymbol{\theta}_1) \cdot p (\boldsymbol{\theta}_1 \mid \boldsymbol{\theta}_2) \cdot p (\boldsymbol{\theta}_2) .$$
⁽²⁵⁾

Our goal is to infer the marginal *a posteriori* distribution $\pi^*(\boldsymbol{\theta}_2)$ of the top level hyperparameters⁴. The marginalization to obtain samples from $\pi^*(\boldsymbol{\theta}_2)$ can

⁴ In the case of the logarithm of the frequency distribution $\log(\nu)$ we are interested in the means and standard deviations of each Gaussian component, but obtaining these parameters from those in Equation 11 by deterministic relationships is straightforward.

be accomplished by a general MCMC procedure in which, once a sample for the joint posterior has been obtained, the procedure retains only the values of θ_2 and discards the rest. The joint posterior distribution of Equation 25 can be efficiently sampled by means of a Gibbs sampling scheme (see Sec. 4.2 of [9]). To reduce our model to the programming language level we have used the BUGS [10] probabilistic language and the OpenBUGS software environment.

3 Application to the *Gaia* Simulated Database of Classical Cepheids

In this Section we evaluate the effectiveness of our model to infer the real distributions of hyperparameters in an extreme scenario of systematic biases in the recovered data. In order to do so, we have constructed a dataset $\mathcal{T} = \{(A_{\text{rec},i}, \nu_{\text{rec},i}, m_{G,\text{rec},i})\}_1^{854} \subseteq \mathcal{D}$ composed of 500 randomly selected instances from the *locus* $\nu_{\text{rec}} = \nu_{\text{in}}$ and all instances (354) from the *locus* $\nu_{\text{rec}} = \pm \nu_{\text{in}} + 7\nu_S - 3\nu_p$. Figure 2 shows the systematic biases for the empirical frequency distribution (histogram) vs the true probability density function (PDF) and for the empirical conditional distributions of the recovered amplitude given the input amplitude for the three *loci* (the identity locus and the $\nu_{\text{rec}} = \pm \nu_{\text{in}} + 7\nu_S - 3\nu_p \ loci$), whose observed parameters are included in the training set.



Fig. 2. Biases in the frequencies (left) and amplitudes (right) present in the training set.

We have trained the model using the OpenBUGS MCMC engine. We have divided the training in two stages and generated three Markov chains (more properly realizations) in each, with a total of 30000 iterations. We have used the first 20000 iterations as a *burn-in* phase, and discarded them after using them for convergence assessment. Thereafter, we obtain 10000 samples from each chain in a second stage (30000 in total). We will assume that these samples were drawn from the posterior distribution of the parameters of interest.

θ	ACR	GRB	$\overline{ heta}$	$2.5\%\mathchar`-97.5\%$ Perc.	Real value
$w_{\nu 2}$	0.39	1.09	0.03	0.01, 0.05	-
$w_{\nu 1}$	0.18	1.03	0.41	0.32, 0.50	-
$w_{\nu 3}$	0.16	1.01	0.57	0.47, 0.66	-
$\mu_{\nu 2}$	0.60	1.16	-1.50	-1.61, -1.37	-
$\mu_{\nu 1}$	0.19	1.01	-0.66	-0.71, -0.61	-
$\mu_{\nu 3}$	0.05	1.01	-0.53	-0.54, -0.51	-
$\sigma_{\nu 2}$	0.83	1.25	0.14	0.10,0.20	-
$\sigma_{\nu 1}$	0.25	1.01	0.28	0.25, 0.33	-
$\sigma_{\nu 3}$	0.09	1.02	0.09	0.08, 0.10	-
a_A	0.04	1.00	-0.43	-0.67, -0.21	-0.5
b_A	0.03	1.00	0.28	-0.02, 0.57	0.2
μ_A	0.00	1.00	0.62	0.58, 0.66	0.7
σ_A	0.00	1.00	0.15	0.14, 0.16	0.15
a_{G1}	0.25	1.04	2.55	2.22, 2.91	-
b_{G1}	0.23	1.03	16.76	16.38, 17.17	-
a_{G2}	0.01	1.00	3.01	2.96, 3.06	-
b_{G2}	0.01	1.00	17.16	17.13, 17.19	-
σ_G	0.00	1.00	0.10	0.09, 0.11	-
λ_{02}	0.00	1.00	-1.132	-1.314 ,-0.955	-
λ_{03}	0.00	1.00	-0.981	-1.156 ,-0.816	-
$\lambda_{\beta 2}$	0.00	1.00	-0.766	-1.140 ,-0.395	-
$\lambda_{\beta 3}$	0.00	1.00	-0.743	-1.091 ,-0.385	-

 Table 2. Summary Statistics of Parameters of Interests.

3.1 Convergence analysis

To evaluate the convergence within and between the three chains we have selected the first 20000 iterations of the algorithm and computed the mean autocorrelation (ACR) (after 200 lags) and the upper bound of a credible interval (at 95%) for the corrected GR statistic [3]. The results of the analysis are summarized in the second and third columns of Table 2. Since the ACR function should decrease to zero as the lag increases and the upper bound for the corrected scale reduction factor (CSRF) should approach unity if the chain is reaching its stationary distribution, we conclude that the worst scenario (high autocorrelation) is encountered in the chains of the parameters specifying the second Gaussian component of log (ν), namely the mixing proportion $w_{\nu 2}$, the mean $\mu_{\nu 2}$ and the standard deviation $\sigma_{\nu 2}$. In particular, chains for $\sigma_{\nu 2}$ show the worst behaviour with a mean ACR after 200 lags of about 0.8 and a CSRF upper bound of 1.25. In contrast, the best scenario is found in the chains of the parameters of the conditional distributions of apparent G-magnitude and amplitude (given the frequency) when $\log(\nu) > -1$, and by chains of logit coefficients. For the slope a_{G2} , the intercept b_{G2} , the mean μ_A and the logit coefficients $\lambda_{\beta j}, \lambda_{0j}, j \in \{1, 2\}$ the mean ACR is nearly zero after lags greater than 50 and the CSRF bound is close to unity.

3.2 Posterior Distributions and Comparison with Real Parameters

In this Section we evaluate the ability of our model to retrieve the real distributions of the frequency, amplitude and apparent G-magnitude of the simulated Cepheids sample from the recovered values in the training set \mathcal{T} . We first compute summary statistics (means and 2.5%-97.5% percentiles) for the samples of the posterior distributions of the hyperparameters inferred by the model. Then, we have compared the posterior means with the parameters of the real theoretical distributions used to generate the simulated sample. Finally, we have constructed theoretical distributions using the posterior means and compared them with the true theoretical distributions and the empirical distribution in the set $\mathcal{I} = \{(A_{\text{in},i}, \nu_{\text{in},i}, m_{G,\text{in},i})\}_{1}^{854}$.



Fig. 3. Posterior versus Real Distributions.

The results of our analysis are shown in Table 2 and Figure 3. We do not include in the Table the parameters used to generate the real frequency $\log (\nu)$, because it is difficult to make a correspondence with the inferred parameters due to the different number of Gaussian components. But if we observe the comparison graph to the left of Figure 3, we conclude that the fitting of $\log (\nu)$ with three components (dotted line), reconstructs the real PDF (solid line) successfully.

For the parameters of the conditional distribution $A_{in} \mid \log(\nu_{in})$ we fitted the piecewise linear model of Equation 14. The middle rows of Table 2 and the graph at the right of Figure 3 show that the system underestimates the true value of the mean μ_A when $\log(\nu_{in}) > -1$.

4 Summary and Conclusions

We have presented a two-level BGM to infer the real distributions of amplitude, frequency and apparent G-magnitude of the Large Magellanic Cloud population of Classical Cepheids from the values recovered by the *Gaia* DPAC pipeline. We have modelled the real frequency by a mixture of three Gaussian distributions and used piecewise linear models (with a fixed knot value depending on the frequency) to model the dependency of the true amplitude and G-magnitude on the true frequency. We have tackled the problem of aliasing in the DPAC frequency recovery module which arises as a result of the Gaia scanning law. We have modelled the recovery probabilities in various *loci* of aliased frequencies using a logistic regression submodel based on the ecliptic latitude predictor. We have modelled the recovered frequencies and amplitudes as generated from mixtures of distributions where the mixing proportions are the recovery probabilities. Although our model has not yet solved completely the aliasing problem (we have only used some predefined configurations of aliased data, and we have restricted the application to a very narrow range of ecliptic latitudes in which the relationship between the recovery probability of aliased frequencies and the ecliptic latitude is monotone) it represents a major step forward. The next step will necessarily consist in extending the analysis to the full celestial sphere by clustering the full variety of time samplings (and corresponding window functions) into discrete bands of ecliptic longitudes and latitudes.

References

- 1. Antonello, E., Fugazza, D., Mantegazza, L.: Variable stars in nearby galaxies. vi. frequency-period distribution of cepheids in ic 1613 and other galaxies of the local group. Astronomy and Astrophysics 388, 477–482 (2002)
- 2. Azzalini, A., Genton, M.G.: Robust likelihood methods based on the skew-t and related distributions. International Statistical Review 76(1), 106–129 (2008)
- 3. Brooks, S.P., Gelman, A.: General methods for monitoring convergence of iterative simulations. Journal of computational and graphical statistics 7(4), 434–455 (1998)
- Eyer, L., Mignard, F.: Rate of correct detection of periodic signal with the gaia satellite. Monthly Notices of the Royal Astronomical Society 361(4), 1136–1144 (2005)
- Gaia Collaboration, Prusti, T., de Bruijne, J.H.J., Brown, A.G.A., Vallenari, A., Babusiaux, C., Bailer-Jones, C.A.L., Bastian, U., Biermann, M., Evans, D.W., et al.: The Gaia mission. Astronomy and Astrophysics 595, A1 (Nov 2016)
- Gelman, A., Jakulin, A., Pittau, M.G., Su, Y.S.: A weakly informative default prior distribution for logistic and other regression models. The Annals of Applied Statistics pp. 1360–1383 (2008)
- 7. Gelman, A., Shalizi, C.R.: Philosophy and the practice of bayesian statistics. British Journal of Mathematical and Statistical Psychology (2013)

- 8. Lauritzen, S.: Graphical Models. Oxford University Press (1996)
- Lunn, D., Jackson, C., Best, N., Thomas, A., Spiegelhalter, D.: The BUGS Book: A Practical Introduction to Bayesian Analysis. CRC Texts in Statistical Science, Chapman & Hall (2012)
- Lunn, D., Spiegelhalter, D., Thomas, A., Best, N.: The bugs project: Evolution, critique and future directions. Statistics in medicine 28(25), 3049–3067 (2009)
- 11. Pearl, J.: Probabilistic Reasoning in Intelligent Systems: Networks of Plausble Inference. Morgan Kaufmann Pub (1988)
- Robert, C.P., Mengersen, K.L.: Reparameterisation issues in mixture modelling and their bearing on mcmc algorithms. Computational Statistics & Data Analysis 29(3), 325–343 (1999)
- 13. Robert, C., Casella, G.: Monte Carlo Statistical Methods. Springer, New York, 2 edn. (2004)
- Sandage, A., Tammann, G., Reindl, B.: New period-luminosity and period-color relations of classical cepheids ii. cepheids in lmc. Astronomy and Astrophysics 424, 43–71 (2004)
- Soszynski, I., Poleski, R., Udalski, A., Szymanski, M.K., Kubiak, M., Pietrzynski, G., Wyrzykowski, L., Szewczyk, O., Ulaczyk, K.: The optical gravitational lensing experiment. the ogle-iii catalog of variable stars. i. classical cepheids in the large magellanic cloud. Acta Astronomica 58, 163–185 (2008)