

# Correcting for sample selection in Bayesian analyses

Samuel R. Hinton,<sup>1,2\*</sup> Alex Kim,<sup>3</sup> Tamara M. Davis,<sup>1,2</sup>

<sup>1</sup>*School of Mathematics and Physics, The University of Queensland, Brisbane, QLD 4072, Australia*

<sup>2</sup>*ARC Centre of Excellence for All-sky Astrophysics (CAASTRO)*

<sup>3</sup>*Lawrence Berkeley National Labs*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

Lifetimes are strictly positive. We cut our data to remove background noise. Astronomical observations suffer from Malmquist bias - our samples are preferentially biased toward those sources that are more likely to be observationally detected. These sample selections should be accounted for. In this paper we present a simple overview of a Bayesian consideration of sample selection, giving a solution to both analytically tractable and intractable models. This can be accomplished via a combination of analytic approximations and Monte Carlo integration, in which dataset simulation is efficiently used to correct for issues in the observed dataset. Toy models are included, along with numerical considerations and optimisations for implementation.

## 1 INTRODUCTION

Sample selection is a problem in many areas of scientific inquiry. It is one of the primary difficulties when performing supernovae cosmology analysis, as our telescopes have visual limits that modify our observed supernovae distribution from the actual underlying distribution. This bias, termed Malmquist bias, is source of much investigation (Butkevich et al. 2005). It is considered during analysis by either modifying the observed data to remove the expected bias (Beto et al. 2014; Conley et al. 2011), or by incorporating the expected bias into the underlying model (Rubin et al. 2015). Truncated data is also commonly encountered in biological fields, where data such as mortality rates are left-truncated (Colchero & Clark 2012). Simplified and generalised examples have been investigated in numerous fashions (Woodroffe 1985; Gull 1989; Grogger & Carson 1991; O’Neill & Barry 1995) and with different fitting algorithms (Gelfand et al. 1992). Whilst generalised resources exist that provide a comprehensive overview of sample selection and analysis techniques in a similar fashion to this work (Klein & Moeschberger 2005), these sources are often opaque due to volume and mathematical complexity.

This work provides a simple treatment of sample selection in a common Bayesian technique. The general theory for considering selection effects is discussed in Section 2. Section 3 provides three examples of increasing complexity with sample selection. Section 4 details numeric concerns and tricks to be aware of for effective implementation of Monte Carlo corrections applied to analytic approximations.

## 2 THEORY

When formulating and fitting a model using a constraining dataset, we wish to resolve the posterior surface defined by

$$P(\theta|\text{data}) \propto P(\text{data}|\theta)P(\theta). \quad (1)$$

Of primary interest to us is the likelihood of observing the data given our parametrised model,  $\mathcal{L} \equiv P(\text{data}|\theta)$ . When dealing with experiments which have selection efficiency, our likelihood necessarily includes that efficiency, for we want to describe the probability that our observations were both drawn from the underlying theoretical model *and* that those observations, given they happened, were subsequently successfully observed. To make this extra conditional explicit, we can write our likelihood as with selection effect  $S$ :

$$\mathcal{L} = P(\text{data}|\theta, S). \quad (2)$$

As we wish to describe our selection efficiency generally as a function dependent on both our data and our model, we can reformulate to

$$\mathcal{L} = \frac{P(S|\text{data}, \theta)P(\text{data}|\theta)}{P(S|\theta)}. \quad (3)$$

Introducing an integral over all possible data to make  $P(S|\theta)$  physical,

$$\mathcal{L} = \frac{P(S|\text{data}, \theta)P(\text{data}|\theta)}{\int P(S, D|\theta) dD} \quad (4)$$

$$= \frac{P(S|\text{data}, \theta)P(\text{data}|\theta)}{\int P(S|D, \theta)P(D|\theta) dD}, \quad (5)$$

where the integral in the denominator has the same dimensionality is the experimental data. **Equation 5 is the generalised likelihood of experiments with sample selection.**

### 3 SAMPLE SELECTION

In this section we present sample selections of increasing complexity. Each sample selection is accompanied by an illustrative example inspired by Type Ia supernova cosmology, where we characterize the properties of a standard or standardizable candle.

#### 3.1 Complete Selection

In a perfect world, data is neither biased nor truncated. The data is perfect. Uncertainties are well quantified and normally distributed around true values. Presumably everything is also spherical and in a vacuum. We thus begin by considering an ideal situation where the sample is complete. All events are included in the sample such that per object,  $P(S|\text{data}, \theta) = 1$ . Trivially this expression is independent of  $\theta$ , and our likelihood from equation (5) reduces down to equation (2). As a concrete example of this case, let us consider a model for a population of objects whose brightnesses form a normal distribution with average  $\mu$  and standard deviation  $\sigma$ . Let us also assume that our experiment produces data  $x$ , measurements of the brightnesses with negligible measurement uncertainty, which we can formally state as

$$\vec{x} \sim \mathcal{N}(\mu, \sigma). \quad (6)$$

If, having collected our observations  $\vec{x}$ , we wanted to constrain  $\mu$  and  $\sigma$ , this would be a simple task of modelling the posterior surface. Taking uniform priors on both parameters, we simply wish to map the surface

$$P(\theta|\text{data}) \propto P(\text{data}|\theta)P(\theta), \quad (7)$$

where our model parameters  $\theta = \{\mu, \sigma\}$  and our data is given by  $\vec{x}$ .

$$P(\mu, \sigma|\vec{x}) \propto P(\vec{x}|\mu, \sigma)P(\mu, \sigma) \quad (8)$$

With uniform priors,  $P(\mu, \sigma) = \text{constant}$ , and can be absorbed into the constant of proportionality. Expanding our observation vector, the posterior surface is given by

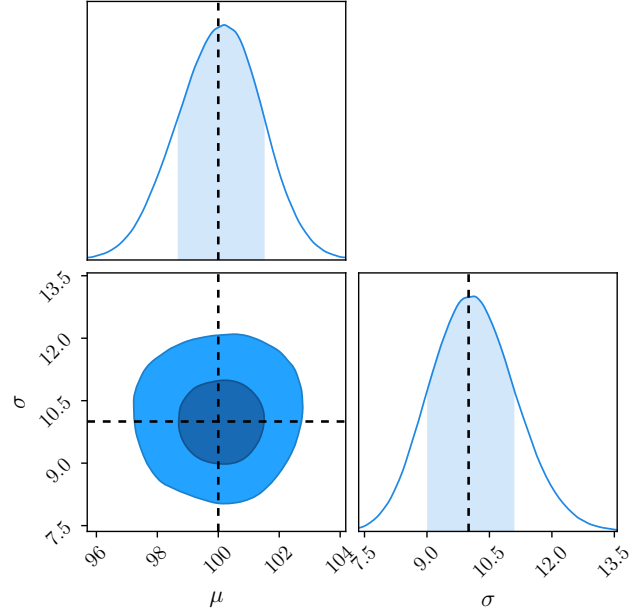
$$P(\mu, \sigma|\vec{x}) \propto \prod_{i=1}^N \mathcal{N}(x_i|\mu, \sigma). \quad (9)$$

Generating a hundred data points with  $\mu = 100$ ,  $\sigma = 10$ , we can recover our input parameters easily, as shown in Figure 1.

#### 3.2 Analytic Sample Selection

We now consider a case with selection bias and has an analytic expression for the likelihood. This illustrative case is useful for the reader as the influence of selection bias is simple and intuitive. The case we consider is identical to the one in the previous section, *except* that only the subset of objects brighter than a threshold  $\alpha$  can be observed.

With this sample selection, all events satisfy  $x > \alpha$ , giving  $P(S|x, \theta) = \mathcal{H}(x - \alpha)$ . We assign a value  $\alpha = 85$  for convenience. If we do not take this truncation into account, we will recover biased parameter estimates. However, we can correct for this truncation using equation (5), as the



**Figure 1.** A systematic test of our perfect model, done by stacking the output chains from fitting 100 independent realisations of our 100 data points. Any systematic offset in our model would be revealed by a shift in the stacked results away from the true parameter values.

integral in the denominator has an analytic solution. Having successfully observed  $x_i$ , it follows that  $x_i > \alpha$  and so  $P(S|x_i, \theta) = 1$ . To substitute in our normal model,

$$\mathcal{L}_i = \frac{P(S|x_i, \theta)P(x_i|\theta)}{\int P(S|D, \theta)P(D|\theta) dD} \quad (10)$$

$$= \frac{\mathcal{N}(x_i|\mu, \sigma)}{\int_{-\infty}^{\infty} \mathcal{H}(D - \alpha) \mathcal{N}(D|\mu, \sigma) dD} \quad (11)$$

$$= \frac{\mathcal{N}(x_i|\mu, \sigma)}{\int_{\alpha}^{\infty} \mathcal{N}(D|\mu, \sigma) dD} \quad (12)$$

$$= \frac{\mathcal{N}(x_i|\mu, \sigma)}{\frac{1}{2} \text{erfc} \left[ \frac{\alpha - \mu}{\sqrt{2}\sigma} \right]}, \quad (13)$$

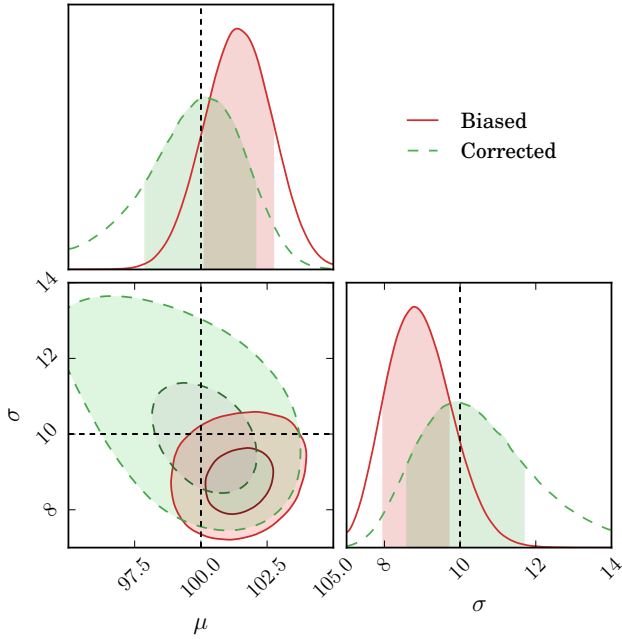
where in the last line we have evaluated the integral in the case  $\mu > \alpha$ . Note that this is for a single observation, and so for a set of independent observations we need to introduce the product found in equation (9).

$$\mathcal{L} = \prod_{i=1}^N \mathcal{L}_i = \prod_{i=1}^N \frac{\mathcal{N}(x_i|\mu, \sigma)}{\frac{1}{2} \text{erfc} \left[ \frac{\alpha - \mu}{\sqrt{2}\sigma} \right]}, \quad (14)$$

However, as our selection efficiency correction is observationally independent it is identical for all observations, allowing us to take it outside the product.

$$\mathcal{L} = 2 \left( \text{erfc} \left[ \frac{\alpha - \mu}{\sqrt{2}\sigma} \right] \right)^{-N} \prod_{i=1}^N \mathcal{N}(x_i|\mu, \sigma), \quad (15)$$

We can add this correction to our model, and note that we now recover unbiased parameter estimates. This is demonstrated in Figure 2, which shows the posterior surfaces for when you take sample selection into account and when you do not. The selection bias preferentially selects intrinsically



**Figure 2.** A systematic test of our imperfect model, done by stacking the output chains from fitting 100 independent realisations of our 100 data points, subject to our thresholding. The bias shown in the red ‘Biased’ contour can be corrected to via the techniques shown in Section 3.2 to recover unbiased surfaces.

brighter objects and, by cutting out some of the distribution, narrows the observed distribution, and so we note that the bias correction correctly increases the weight of low  $\mu$  and high  $\sigma$  parametrisations, as those models would be subject to the most sample bias. We also note that not only does the best fit location for each parameter shift, but the shape of the posterior surface itself is significantly modified.

### 3.3 Analytically Intractable Sample Selection

Unfortunately it is a rare scenario when dealing with nature and all her faults for us to have an analytic selection function, let alone a function encapsulated by a single parameter. A more realistic scenario involves a selection efficiency instead would take the form of non-analytic function of many model parameters. And the function would probably be stochastic too, just to throw another wrench in the works. Provided a method of forward modelling or simulating observations, the solution is to combine an analytic approximate correction with Monte Carlo integration.

In this subsection we introduce an example that presents computational challenges in the evaluation of the posterior, and which better reflects the model complexity in supernova cosmology analysis. The specific challenge is that the denominator in the equation for the likelihood equation (5) is dependent on  $\theta$  and does not have an analytic solution (i.e. has to be solved numerically). When using Monte Carlo sampling methods to probe the posterior with this kind of likelihood, the integral in the denominator must be calculated with each draw. We therefore present importance sampling as an approach that make these posteriors more computationally tractable.

Let us modify our imperfect toy model from the previous section. Instead of observing just one variable,  $x$ , we also observe a new independent variable,  $y$ , which is drawn from its own distribution  $y \sim \mathcal{N}(\mu_y, \sigma_y)$ , and has no measurement uncertainty (like  $x$ ). Our selection efficiency can now become a combination of  $x$  and  $y$ , such that we only observe events that satisfy  $x + \beta y > \alpha$ , giving  $P(S|x, y, \theta) = \mathcal{H}(x + \beta y - \alpha)$ . Our likelihood for such a toy model becomes now the combination of probabilities for observing both  $x$  and  $y$ , with the denominator becoming an integral over all possible  $X$  and  $Y$  observations subject to our selection effects.

$$\mathcal{L}_i = \frac{\mathcal{N}(x_i|\mu, \sigma)\mathcal{N}(y_i|\mu_y, \sigma_y)}{\int_{-\infty}^{\infty} \mathcal{H}(X + \beta Y - \alpha) \mathcal{N}(X|\mu, \sigma) \mathcal{N}(Y|\mu_y, \sigma_y) dX dY} \quad (16)$$

Assume that we cannot solve this integral analytically, and must resort to numeric solutions. These often clash with sampling methods, especially for high dimensional integrals. Inserting Monte Carlo integration into fitting algorithms can drastically slow them down, and algorithms such as Hamiltonian MCMC that require continuous surfaces can easily fail on surfaces that fluctuate from Monte Carlo integration. Even by fixing the samples used in MC integration (thereby giving a continuous surface), the complexity of the surface derivatives will pose almost insurmountable problems for any algorithms that utilise surface gradients.

One solution is to find an approximate, analytic correction we can utilise in our fitting algorithm which seeks to shift the region of parameter space sampled by the sampler closer to the correct area, and then importance sample our MC chains to provide a fully corrected surface. Given an approximate analytic correction  $w_{\text{approx}}$ , we explicitly break our likelihood into two parts,  $\mathcal{L}_i = \mathcal{L}_{i1}\mathcal{L}_{i2}$ , with the parts given by

$$\mathcal{L}_{i1} = \frac{P(S|\text{data}, \theta)P(\text{data}|\theta)}{w_{\text{approx}}} \quad (17)$$

$$\mathcal{L}_{i2} = \frac{w_{\text{approx}}}{\int P(S, D|\theta) dD} \quad (18)$$

Where  $\mathcal{L}_{i1}$  represents an analytic function which can efficiently be sampled using traditional Monte Carlo sampling methods, and  $\mathcal{L}_{i2}$  is a numerically calculated weight applied after Monte Carlo sampling in order to model the full likelihood surface.

In our example, if  $\beta \ll 1$ , such that the majority of selection effect is encapsulated by  $x$  and not  $y$ , our approximate correction can take the form found in the previous correction from Section 3.2. Having true values of  $\mu = 100$ ,  $\sigma = 10$ ,  $\mu_y = 30$ ,  $\sigma_y = 5$ , and a known  $\beta = 0.2$ , we can give a concrete example. Assuming some prior, imperfect knowledge of  $\mu_y$  (perhaps we believe it is approximately 20) we estimate that the average contribution from  $\beta y$  is approximately  $20\beta = 4$  (which is close to the correct value of 6), and from this we add in a small adjustment to the analytic correction from Section 3.2:

$$w_{\text{approx}} = \frac{1}{2} \text{erfc} \left[ \frac{\alpha - \mu - 4}{\sqrt{2}\sigma} \right]. \quad (19)$$

This gives our likelihood parts as

$$\mathcal{L}_{i1} = \frac{\mathcal{N}(x_i|\mu, \sigma)\mathcal{N}(y_i|\mu_y, \sigma_y)}{w_{\text{approx}}} \quad (20)$$

$$\mathcal{L}_{i2} = \frac{w_{\text{approx}}}{\int \int_{-\infty}^{\infty} \mathcal{H}(x + \beta y - \alpha) \mathcal{N}(X|\mu, \sigma) \mathcal{N}(Y|\mu_y, \sigma_y) dX dY}. \quad (21)$$

As stated previously,  $\mathcal{L}_1$  can thus be fitted with a traditional sampler without numeric difficulty or slowdown, and  $\mathcal{L}_2$  allows us to calculate the weight of each sample. We are effectively importance sampling our likelihood evaluations. The computational benefits of this should not be understated either - each sample in our chains can be reweighted independently, providing a task that is trivially parallelisable. Evaluating  $\mathcal{L}_2$  using Monte Carlo integration of  $n$  samples, we have

$$\mathcal{L}_{i2} = \frac{w_{\text{approx}} n}{\sum_{j=1}^n \mathcal{H}(X_j + \beta Y_j - \alpha) \mathcal{N}(X_j|\mu, \sigma) \mathcal{N}(Y_j|\mu_y, \sigma_y)}. \quad (22)$$

We can now easily move from a single observation to a set of  $N$  observations.

$$\mathcal{L}_1 = w_{\text{approx}}^{-N} \prod_{i=1}^N \mathcal{N}(x_i|\mu, \sigma) \mathcal{N}(y_i|\mu_y, \sigma_y) \quad (23)$$

$$\mathcal{L}_2 = \left( \frac{w_{\text{approx}} n}{\sum_{j=1}^n \mathcal{H}(X_j + \beta Y_j - \alpha) \mathcal{N}(X_j|\mu, \sigma) \mathcal{N}(Y_j|\mu_y, \sigma_y)} \right)^N \quad (24)$$

Thus we end up with a corrected posterior surface as shown in Figure 3. It is important to note that the point of maximum likelihood is biased by roughly the same amount in the biased and approximate posterior surfaces in Figure 3. However, as the approximately correct posterior is broader than the biased surface, it has far more samples in the region of parameter space mapped by the fully corrected posterior. This is the entire purpose of the approximate correction - to maximise the number of samples in the correct region of parameter space, so that we can importance sample our chains efficiently.

## 4 NUMERICAL TRICKS

### 4.1 Importance Sampling

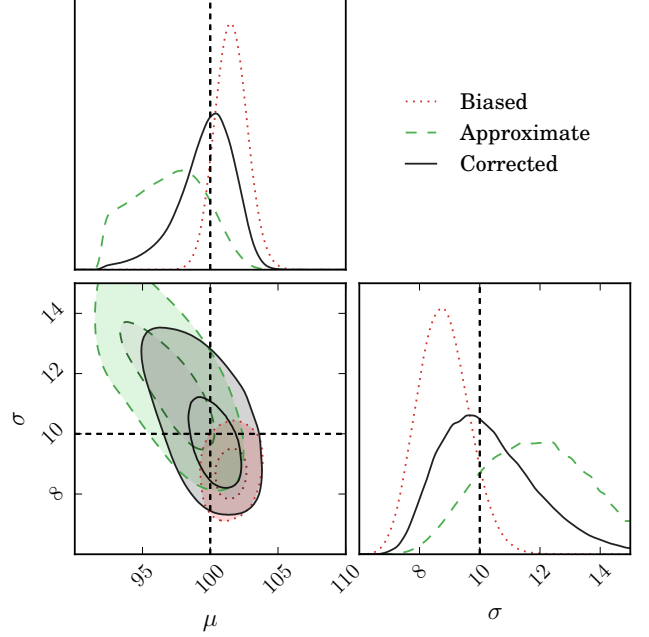
Further tricks can be used to increase the efficiency with which the samples are reweighted. Firstly, the overarching analytic model often provides functions which can be drawn from efficiently. In the case of our example, by drawing the random numbers  $X$  and  $Y$  respectively from the normal distributions  $\mathcal{N}(\mu, \sigma)$  and  $\mathcal{N}(\mu_y, \sigma_y)$  (ie traditional importance sampling) we need only evaluate the step function for our data points. That is, we replace

$$\int f(X) \mathcal{N}(X|\mu, \sigma) dX = \frac{1}{N} \sum_{i=1}^N f(X_i) \mathcal{N}(X_i|\mu, \sigma), \quad (25)$$

where  $X$  is drawn from a uniform distribution over all space, with

$$\int f(X) \mathcal{N}(X|\mu, \sigma) dX = \frac{1}{N} \sum_{i=1}^N f(X_i), \quad (26)$$

where  $X$  drawn from  $\mathcal{N}(X_i|\mu, \sigma)$ .



**Figure 3.** A systematic test of our more complicated model, done by stacking the output chains from fitting 100 independent realisations of our 100 data points, subject to our thresholding. The likelihood  $\mathcal{L}_1$  was evaluated, and reweighted using Monte Carlo integration of a hundred thousand possible events as per  $\mathcal{L}_2$ . The truncated data with no correction is shown as ‘Biased’ in dotted red, the ‘Approximate’ only correction ( $\mathcal{L}_1$ ) shown in dashed green, and the final reweighted chain shown in solid black as ‘Corrected’.

### 4.2 Precomputing selection

If evaluating the probability that an event is observed is numerically expensive (i.e. not a step function), it is easy to pregenerate a set of events and reuse them for all weights - provided that the number of events used when calculating the weights is sufficient to make the statistical error of Monte Carlo integration insignificant when compared to the constraining power of your dataset. This method is however only efficient when prior knowledge of parameter values is known to allow a reasonable initial draw of events. Without this prior information, samples need to span the entire posterior volume, which is numerically intractable even for low dimensional models.

Consider the imperfect example - where we observe  $x$  drawn from an underlying normal distribution, but utilise the Monte Carlo integration technique from Section 3.3 and do not have an analytic approximation (i.e. we set  $w_{\text{approx}} = 1$ ). We could estimate, given some prior knowledge, that variable  $x \approx \mathcal{N}(\mu_{\text{guess}}, \sigma_{\text{guess}})$ . We then draw samples of  $x$  from this distribution, recording the probability of each draw and then calculating whether our potential observation of  $x$  would be observed in the experiment or not. That is, we assign  $P(S|x, \theta) = P(S|x) = 1$  or  $0$  given it passed cuts or not. We discard all events with  $0$  weight (as they have  $0$  weight), and only track those events which pass. Then, when calculating the sample reweighting after running chains,  $\mathcal{L}_{i2}$

becomes

$$\mathcal{L}_{i2} \propto w_{\text{approx}} \left[ \sum_{j=1}^n \frac{\mathcal{N}(X_j|\mu, \sigma)}{\mathcal{N}(X_j|\mu_{\text{guess}}, \sigma_{\text{guess}})} \right]^{-1}, \quad (27)$$

where you can see that we discard the constant  $n$  from equation (22) as we only care about likelihood proportionality. Provided our parameter estimate is reasonably well informed, the computation benefit this precomputation provides is enormous for any nontrivial selection function. Not only do we now waste no time when calculating  $\mathcal{L}_{i2}$  determining  $P(S|\text{data}, \theta)$ , as we only save results that pass the cuts, we have no wasted evaluations of  $\mathcal{N}(X_j|\mu, \sigma)$ .

The astute reader may have picked up on one assumption - that selection efficiency of an observation is independent of model parameters  $\theta$ . For many experimental cases this should hold, however if it does not this method cannot be used to increase efficiency. Gridding or interpolating the parameter space is strongly not recommended due the required accuracy of  $\mathcal{L}_{i2}$ . Even a small error when raised to the power of  $N$  can spiral out of control.

### 4.3 Log-space

Following from the previous section, as our reweighting  $\mathcal{L}_2$  is raised to the power of the number of our observations, they should definitely be computed in log-space, which turns the power into a linear factor. As most probabilistic work is already computed in log-space, this subsection barely needs to be stated. However, whilst working in log-space an efficient way of increasing the accuracy of the approximate analytic correction is to fit the correction such that the spread of the distribution  $\log \mathcal{L}_2$  is minimised.

## 5 CONCLUSION

Sample selection is a pervasive issue in many scientific domains. For simple cases of sample selection which can be encapsulated with analytic functions, it is possible to analytically correct likelihood surfaces by introducing selection efficiency into the likelihood formulation. When analytic corrections fail to provide an adequate description of selection effects, Monte Carlo integration can be used on top of analytic approximations to further correct the likelihood surface, provided the experiment can be effectively simulated.

## ACKNOWLEDGMENTS

We gratefully acknowledge the input of the many researchers that were consulted during the creation of this paper.

## REFERENCES

- Betoule M., et al., 2014, *A&A*, **568**, A22  
 Butkevich A. G., Berdyugin A. V., Teerikorpi P., 2005, *Monthly Notices of the Royal Astronomical Society*, **362**, 321  
 Colchero F., Clark J. S., 2012, *Journal of Animal Ecology*, **81**, 139  
 Conley A., et al., 2011, *ApJS*, **192**, 1  
 Gelfand A. E., Smith A. F. M., Lee T.-M., 1992, *Journal of the American Statistical Association*, **87**, 523

- Grogger J. T., Carson R. T., 1991, *Journal of applied econometrics*, **6**, 225  
 Gull S. F., 1989, in , *Maximum Entropy and Bayesian Methods*. Springer, pp 511–518  
 Klein J. P., Moeschberger M. L., 2005, *Survival analysis: techniques for censored and truncated data*. Springer Science & Business Media  
 O'Neill T. J., Barry S. C., 1995, *Biometrics*, pp 533–541  
 Rubin D., et al., 2015, *ApJ*, **813**, 137  
 Woodroffe M., 1985, *The Annals of Statistics*, pp 163–177

This paper has been typeset from a  $\text{\TeX/L\AA T\TeX}$  file prepared by the author.