# Robust Bayesian regression in astronomy

William Martin[1]⋆ and Daniel Mortlock[1,2]

[1]*Department of Physics, Imperial College London, Blackett Laboratory, Prince Consort Road, London SW7 2AZ, UK*
[2]*Department of Mathematics, Imperial College London, London, SW7 2AZ, UK*

## ABSTRACT

Model mis-specification (e.g. the presence of outliers) is commonly encountered in astronomical analyses, often requiring the use of ad hoc algorithms (e.g. sigma-clipping). We develop and implement a generic Bayesian approach to linear regression, based on Student's $t$-distributions, that is robust to outliers and mis-specification of the noise model. Our method is validated using simulated datasets with various degrees of model mis-specification; the derived constraints are shown to be systematically less biased than those from a similar model using normal distributions. We demonstrate that, for a dataset without outliers, a worst-case inference using $t$-distributions would give unbiased results with $\lesssim 10$ per cent increase in the reported parameter uncertainties. We also compare with existing analyses of real-world datasets, finding qualitatively different results where normal distributions have been used and agreement where more robust methods have been applied. A Python implementation of this model, $t$-cup, is made available for others to use.

**Key words:** methods: statistical – methods: data analysis – software: data analysis

## 1 INTRODUCTION

Linear regression is a common problem in astronomy, arising in fields as diverse as galaxy formation and evolution (e.g. the supermassive black hole (SMBH) mass – stellar velocity dispersion correlation, Ferrarese & Merritt 2000; Gebhardt et al. 2000), stellar physics (e.g. the Leavitt law linking the luminosity and pulsation period for Cepheid variable stars, Leavitt & Pickering 1912), and cosmology (e.g. the original formulation of Hubble's law, Hubble 1929). For this reason, there are a plethora of techniques used by astronomers for linear regression when both the dependent and independent variables are measured with error (e.g., Press et al. 1992; Akritas & Bershady 1996; Tremaine et al. 2002; Kelly 2007 — see Andreon & Hurn 2013 and Andreon & Weaver 2015 for reviews including further examples).

Kelly (2007) illustrates that common ad-hoc estimators such as FITEXY (Press et al. 1992; Tremaine et al. 2002) and BCES (Akritas & Bershady 1996) suffer from biases and can underestimate intrinsic scatter. Similar algorithms exist for removing outliers from data (e.g. sigma clipping) have been used previously but these can lead to controversy (e.g. the reanalysis of the data-set from Riess et al. 2011 by Efstathiou 2014). A more principled approach is to model the entire measurement process.

Bayesian hierarchical models (BHMs) are a natural way to model such datasets — these allow astronomers to account for, e.g., measurement errors, selection effects, interlinked parameters, censored data, and many other effects common to astronomical problems. BHMs have seen increasing use in astrophysics over the past few decades, from the distance-redshift relation in cosmology (e.g. Feeney et al. 2018; Avelino et al. 2019) and photometric redshift estimation (e.g. Leistedt et al. 2016) to exoplanet characterisation (e.g. Sestovic et al. 2018) and population-level inference (e.g. Kelly et al. 2009) — see Feigelson et al. (2021) for a recent review including further examples of BHMs.

Kelly (2007) presented a general BHM for linear regression with measurement errors and censored data, demonstrating several advantages over the other methods considered: no bootstrapping was required to obtain uncertainties on parameters; the Bayesian approach was easily extensible to truncated or censored data; and other methods would sometimes severely underestimate intrinsic scatter in the data. This formulation of Bayesian regression, sometimes known as LINMIX_ERR, is now commonly used in astronomy (e.g. McConnell & Ma 2013; Bentz et al. 2013; Andrews et al. 2013). This approach has been extended and refined by others (e.g. Mantz 2016; Sereno 2016; Bartlett & Desmond 2023; Jing & Li 2024). These models assume parameters are normally distributed throughout; however, scientific uncertainties are often empirically not normally distributed, leading to more frequent outliers (Bailey 2017). Inference that relies on normal distributions can be unduly affected by outliers (see Section 4 for an exploration of this effect).

The problem of outliers within datasets can be thought of as model mis-specification: these objects do not fit the distributions used to model them. The ideas behind robust inference can prove useful for this problem; in robust inference, methods are designed to work irrespective of the actual generative distribution (Berger et al. 1994). One approach is to use distributions that are leptokurtic (i.e. have heavier tails than a normal distribution) for inference — examples include Student's $t$-distributions (Andrews & Mallows 1974) or Gaussian mixture models (Box & Tiao 1968; Aitkin & Tunnicliffe Wilson 1980). Leptokurtic distributions can lead to more robust results in the case of model misspecification (e.g. Berger et al. 1994; Sivia & Skilling 2006; Gelman et al. 2013). Student's $t$-distributions in

⋆ E-mail: w.martin19@imperial.ac.uk

particular have seen use in bespoke astronomical (e.g. Andreon et al. 2008; Jontof-Hutter et al. 2016; Park et al. 2017; Andreon 2020) and cosmological (e.g. Feeney et al. 2018) inference, but there is not currently a generic robust method for Bayesian astronomical data analysis.

We propose a development of a generic approach for robust astronomical data analysis. From the review of previous methods, we can identify properties that we would like to see in our regression model:

• A BHM – we favour a hierarchical approach because it naturally encodes the hierarchical structure of astronomical regression problems (i.e. objects are drawn from a high-level population; the objects intrinsically obey some relationship; the objects are measured with error and only the measured values are known). We adopt a Bayesian approach both for practical reasons – the resultant posterior distributions provide full uncertainty quantification – and due to their logical consistency (e.g. Cox 1946; Van Horn 2003; Knuth & Skilling 2010).

• A robust model – we desire a model that is robust to both outliers (i.e. data points that, whether intrinsically or by virtue of measurement errors, do not lie sufficiently close to our regression relation) and model mis-specification (i.e. where the underlying distribution of data does not match up with the distribution assumed for modelling).

• A general method – we seek a model that does not require case-by-case optimisation for application to different regression problems (e.g. no manual outlier identification and removal, no need to rescale prior distributions for different problems, etc.).

We implement these ideas in this paper, beginning with a discussion of our model in Section 2. In Section 3, we outline the methods that we use to validate our model; the results of these validation checks are presented in Section 4. In Section 5, we compare the performance of our model on real-world datasets with the models outlined in Kelly (2007) and Park et al. (2017), before summarizing our conclusions in Section 6.

# 2 FORMALISM

Here we establish notation and set out our model. Our dataset has $N$ astronomical objects, each with a $K$-dimensional vector of associated independent quantities $\{x_i\}$ and a dependent quantity $\{y_i\}$. For example, if we were estimating SMBH mass using measurements of luminosity and line width, we would have $K = 2$ dimensional vectors of independent quantities $\{x_i\} = \{(L_i, \Delta V_i)^T\}$, and the dependent quantity $\{y_i\}$ would correspond to the black hole mass.

We assume that the independent quantities $\{x_i\}$ and dependent quantities $\{y_i\}$ are related by a regression relation

$$y_i = f(x_i; \theta_f) + \delta_i, \tag{1}$$

$$\delta_i \sim \mathcal{P}_{\text{int}}(\theta_{\text{int}}), \tag{2}$$

where $f(x_i; \theta_f)$ is a function relating $\{x_i\}$ to $\{y_i\}$, with parameters $\theta_f$, and $\mathcal{P}_{\text{int}}$ is an unknown probability distribution with parameters $\theta_{\text{int}}$.

These objects are then observed, resulting in the measured data

$$\hat{x}_i = x_i + \epsilon_{x,i}, \tag{3}$$

$$\hat{y}_i = y_i + \epsilon_{y,i}, \tag{4}$$

$$\epsilon_{x,i} \sim \mathcal{P}_{\text{obs}}(\theta_{\text{obs}}), \tag{5}$$

$$\epsilon_{y,i} \sim \mathcal{P}_{\text{obs}}(\theta_{\text{obs}}), \tag{6}$$

where $\mathcal{P}_{\text{obs}}$ is an unknown probability distribution with parameters $\theta_{\text{obs}}$.

In a Bayesian framework, we can extend this model to deal with,

e.g., censored data or selection effects, but this is beyond the scope of the current work — see Kelly (2007) for an overview of an approach that would incorporate these effects.

## 2.1 Building a robust model

In the setup outlined above, the form of the distributions $\mathcal{P}_{\text{int}}$ and $\mathcal{P}_{\text{obs}}$ are unknown. A common assumption in analysis is to use normal distributions to model both of these distributions. However, in the case of model mis-specification (where, e.g., we assume $\mathcal{P}_{\text{int}}$ is a normal distribution but it is a different distribution), the results obtained under this assumption can be biased. Making the assumption that $\mathcal{P}_{\text{int}}$ follows a different distribution can give results that are robust to this model misspecification. This means that, even though we do not believe that the distribution we choose to model $\mathcal{P}_{\text{int}}$ is the same as the underlying, unknown $\mathcal{P}_{\text{int}}$, we can be confident in the inferences about the regression model and the properties of individual objects.

In this paper, we use Student's $t$-distributions to build a model that is robust to model misspecification.

## 2.2 Sampling distribution

For robust inference, we seek a sampling distribution that can have heavier tails than a normal distribution, but that can reduce to a normal distribution when the underlying dataset is normally distributed. We further seek an identifiable model, and a differentiable distribution which can be fit using Hamiltonian Monte Carlo (HMC). These constraints lead us naturally to Student's $t$-distributions, which fulfil all of these criteria.

Student's $t$-distributions are encountered when estimating the mean of a normal distribution with unknown variance from a limited number of samples. The number of samples is a parameter of the distribution: for $n$ samples, the corresponding Student's $t$-distribution will have $\nu \equiv n - 1$ "degrees-of-freedom". This value of $\nu$ parameterizes how heavy-tailed the distribution is. While the interpretation of $\nu$ as "degrees-of-freedom" only makes sense for $\nu \in \mathbb{Z}^+$, the distribution is normalizable for any positive, real $\nu$; for this reason, we shall refer to $\nu$ as the shape parameter.

The Student's $t$-distribution, with location $\mu$ and scale $\sigma$, has the probability density function

$$t_\nu\left(x; \mu, \sigma^2\right) = \frac{1}{\sqrt{\pi\nu}\sigma} \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{1}{\nu} \frac{(x-\mu)^2}{\sigma^2}\right)^{-\frac{\nu+1}{2}}. \tag{7}$$

This is shown for a range of $\nu$ in Figure 1: $\nu = 1$ gives a Cauchy distribution; and $\nu \to \infty$ tends to a normal distribution. The distribution has mean

$$\mathbb{E}(x) = \begin{cases} \mu & \nu > 1, \\ \text{undefined} & \text{otherwise} \end{cases} \tag{8}$$

and variance

$$\text{Var}(x) = \begin{cases} \frac{\nu}{\nu-2}\sigma^2 & \nu > 1, \\ \infty & 1 < \nu \le 2, \\ \text{undefined} & \text{otherwise}. \end{cases} \tag{9}$$

In this paper, we have found it useful to define two quantities for comparison with normal distributions. The first, $\sigma_{68}(\nu)$, is the width of the highest density interval for a $t$-distribution with scale parameter $\sigma = 1$ such that the density contained in the interval is
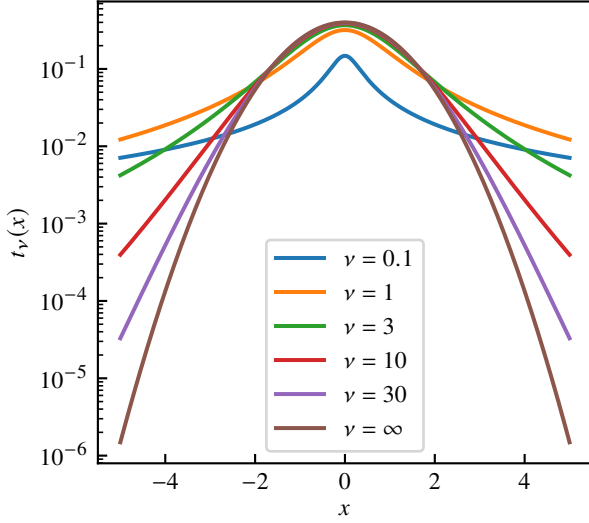
**Figure 1.** The $t$-distribution probability density function for different values of shape parameter, $\nu$.

equal to that of a $1\sigma$ interval for a normal distribution. This is given by

$$\sigma_{68}(\nu) = \sqrt{\nu \left( \frac{1}{I^{-1}\left(\Phi_{\sigma=1}; \frac{\nu}{2}, \frac{1}{2}\right)} - 1 \right)}, \qquad (10)$$

where $\Phi_{\sigma=1} = \text{erf}\left(\frac{\sqrt{2}}{2}\right) \approx 0.6827$ is the posterior density contained within $\pm 1\sigma$ of the mean for a normal distribution, and $I^{-1}$ is the inverse of the regularized incomplete beta function, which we define as

$$I(x; a, b) = \frac{\int_0^x t^{a-1}(1-t)^{b-1} \mathrm{d}t}{\int_0^1 t^{a-1}(1-t)^{b-1} \mathrm{d}t}. \qquad (11)$$

This quantity is useful for defining a "scale" parameter that is less tightly coupled to $\nu$ than the $\sigma$ parameter of the $t$-distribution.

The second quantity we use is an "outlier fraction" $\omega$, which is defined as the fraction of points expected to lie more than $3\sigma$ from the distribution mean $\mu$. For a $t$-distribution, this is given by

$$\omega(\nu) = P\left(|x - \mu| > 3\sigma \mid x \sim t_\nu(\mu, \sigma)\right) \qquad (12)$$
$$= 2F_\nu(\mu - 3\sigma), \qquad (13)$$

where $F_\nu(x)$ is the cumulative distribution function for the Student's $t$-distribution with shape parameter $\nu$. Figure 2 shows the relationship between outlier fraction $\omega$ and shape parameter $\nu$, with $\omega \to 2.70 \times 10^{-3}$ in the limit of a normal distribution, i.e., as $\nu \to \infty$. Under this definition approximately 1 in 370 data-points would be outliers for data that follows a normal distribution; for Cauchy-distributed data (i.e. $\nu = 1$), every fifth data-point to be an outlier.

## 2.3 Regression model

Our regression model, represented as a directed acyclic graph in Figure 3, is specified here. We assume a linear relationship between $\{x_i\}$ and $\{y_i\}$, i.e.

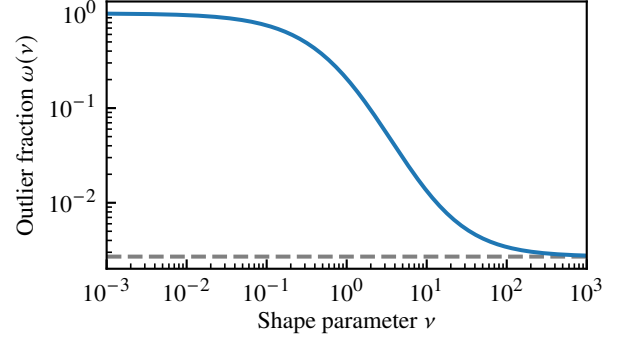$$f(x_i; \theta_f) = \alpha + \beta \cdot x_i \qquad (14)$$



**Figure 2.** The outlier fraction, $\omega$, for different values of shape parameter, $\nu$. The outlier fraction of a normal distribution ($\approx 2.70 \times 10^{-3}$) is indicated by the dashed grey line.

where $\{\alpha, \beta\} \equiv \theta_f$. In this paper, we only consider this linear relationship, but this apparatus could be used for other, non-linear relationships between variables with different parameters.

Similarly to Kelly (2007), we define a Bayesian hierarchical model to reflect the nature of the regression problem. Firstly, we assume that we can represent the probability distribution $\mathcal{P}_{\text{int}}$ with a Student's $t$-distribution with shape parameter $\nu$ and scale parameter $\sigma_{\text{int}}$ — i.e.

$$y_i \sim t_\nu\left(\alpha + \beta x_i, \sigma_{\text{int}}^2\right). \qquad (15)$$

We use a Student's $t$-distribution here not because we believe that the data intrinsically follows this distribution, but because the resulting model is robust to model mis-specification — such as if a galaxy were to be an outlier from a particular relation as a result of a recent merger.

We assume that we can represent the probability distribution $\mathcal{P}_{\text{obs}}$ as a normal distribution with scale parameter given by the error associated with the measured quantity. These measurements are modelled as

$$\hat{x}_i \sim \mathcal{N}\left(x_i, \Sigma_{x,i}^2\right) \qquad (16)$$
$$\hat{y}_i \sim \mathcal{N}\left(y_i, \sigma_{y,i}^2\right). \qquad (17)$$

We experimented with assuming $\mathcal{P}_{\text{obs}}$ to be $t$-distributed, but found that the resultant model was difficult to sample as a result of its geometry, and that posterior predictive checks often included large outliers and did not resemble the datasets that gave rise to them. Fortunately, for robust inference the presence of a single heavy-tailed component is sufficient.

## 2.4 Priors

To ensure that our model is generically applicable to astronomical linear regression, we ensure that both independent and dependent quantities are scaled to have zero mean and unit variance. This allows us to set generic priors on the scaled intercept, gradients and scatter $\{\tilde{\alpha}, \tilde{\beta}, \tilde{\sigma}_{\text{int}}\}$ that do not require rescaling for different units or datasets (though these priors can be revised to incorporate prior information).

Our default priors on the regression parameters $\{\tilde{\alpha}, \tilde{\beta}, \tilde{\sigma}_{\text{int}}\}$ are

$$\tilde{\alpha} \sim \mathcal{N}(\mu = 0, \sigma^2 = 4), \qquad (18)$$
$$\tilde{\beta} \sim \mathcal{N}(\mu = 0, \sigma^2 = 4), \qquad (19)$$
$$\tilde{\sigma}_{68} \sim \Gamma(1.1, 5). \qquad (20)$$

**Figure 3.** A directed acyclic graph representing the $t$-cup model for the function $f(\mathbf{x};\boldsymbol{\theta}) = \alpha + \boldsymbol{\beta} \cdot \mathbf{x}_i$.

The motivation behind the prior on $\tilde{\alpha}$ is that, as the data is pre-scaled to have zero mean and the relationship between quantities is assumed to be linear, we would expect the data to have zero intercept. Similarly, the prior on $\tilde{\beta}$ has been chosen because the pre-scaling of the data suggests a gradient between $-1$ and $1$. The prior on $\tilde{\sigma}_{68}$ is informed by Chung et al. (2013), who demonstrate that a prior with density at $\tilde{\sigma}_{\text{int}} = 0$ will lead to a prediction of no intrinsic scatter in datasets where intrinsic scatter is present; in addition, there is a reasonable physical argument that most astrophysical processes will not produce objects with no population scatter. The chosen prior balances this constraint with difficulties arising in testing as a result of the pre-scaling, which led to insufficient prior density at small levels of intrinsic scatter.

The prior on the latent values of the independent quantities $\{\mathbf{x}_i\}$ is a Gaussian mixture model prior, similar to that used in Kelly (2007). While Kelly infers the complete prior as part of the model, we found that such a prior led to sampling issues as a result of the lack of identifiability inherent in a multicomponent Gaussian mixture model. Therefore, we use extreme deconvolution (Bovy et al. 2011) to estimate the parameters of a Gaussian mixture model approximating the latent distribution of $\{\mathbf{x}_i\}$, and use these parameters for the prior.

The prior on the shape parameter $\nu$ is

$$\nu \sim \text{Inv-}\Gamma(4, 15); \tag{21}$$

this prior was chosen to balance flexibility (such that the model could accommodate both heavily leptokurtic distributions — e.g. the Cauchy distribution — and normally-distributed data) and sampling issues (as $\nu \leftarrow \infty$, different values of $\nu$ become rapidly indistinguishable). The reasoning behind this choice of prior is expanded upon in Appendix A. The exact choice of prior for $\nu$ is unimportant, as the model is designed to be robust to model mis-specification; the only key element is that $\nu$ is allowed to vary to accommodate both normal and heavy-tailed distributions. By virtue of this, difficulties in sampling $\nu$ according to this prior distribution do not invalidate

the results of the model, as long as there is coverage at both small and large values of $\nu$.

### 2.5 Asymptotic normality

One potential disadvantage of adopting a heavy-tailed likelihood in the statistical model is that the resultant inferences would, in general, be less constraining than under the assumption of a normal model (see, e.g., the motivation behind the mixture model proposed in Tak et al. 2019). If there are no outliers, the implication would be that the robust approach is inferior. This effect is illustrated in Figure 4, which shows results for datasets of $N = 100$ points drawn from a normal distribution with zero mean and unit variance. The points show the ratio of the posterior standard deviations for the mean when analyzed using a $t$-distribution to that obtained using a normal distribution. These agree well with the approximate form of this ratio obtained using the Fisher information of a $t$-distribution likelihood (see Appendix B). This ratio increases with decreasing $\nu$ until $nu \sim 0.6$, beyond which it falls again; this is a result of the increasingly narrow peak of the $t$-distribution (see Figure 1). While the uncertainty can be increased by up to $\sim 20\%$ (in the Cauchy regime, $\nu \approx 1$), we consider this an acceptable trade-off to reduce bias in cases of model mis-specification. Under a flexible model that infers $\nu$ (such as the one presented here), we would expect the model to infer a value for $\nu$ that approaches a normal distribution, mitigating this issue of increased error.

This effect would be appreciable for the smallest datasets – the posteriors would have heavier tails than under a normal model – but, in this case, meaningful constraints are unlikely irrespective of the adopted model. For larger datasets, of more than a few tens of points, these tails are effectively multiplied out, typically leaving just a single region of high probability. Moreover, for the regression problem considered here, the posterior in the regression parameters satisfies the requirements for asymptotic normality (e.g. Ghosh et al. 2006), so the core of the posterior is Gaussian in form, just as would be the case under a normal model. (There is not even a significant numerical cost as the posterior has to be evaluated by sampling anyway.)

### 2.6 Implementation

The model described in Section 2.3 is implemented in NumPyro (Phan et al. 2019; Bingham et al. 2019), which is used to draw samples via a HMC No U-Turn Sampler. The implementation is packaged as $t$-cup, available as a Python package[1].

For the purposes of comparison, we have also implemented a model that mirrors the above structure, but uses normal distributions at each stage — we refer to this as $n$-cup.

### 3 VALIDATION

In this section, we outline the methods used to validate the model set out in Section 2. We run two types of tests to validate the performance of $t$-cup under different types of model mis-specification: simulation-based calibration tests (Cook et al. 2006; Talts et al. 2018); and fixed value calibration tests.

---
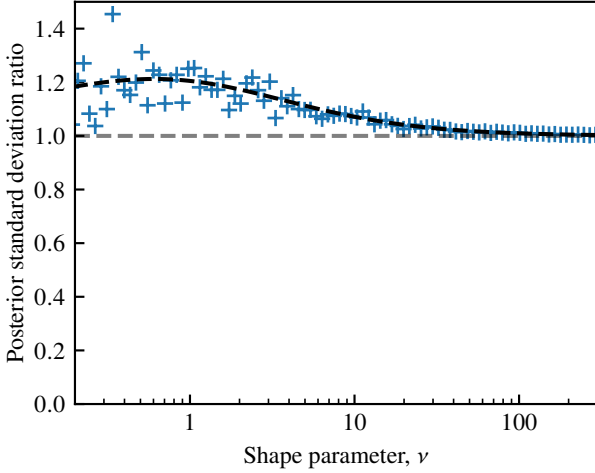
[1] https://github.com/wm1995/tcup

**Figure 4.** The ratio of the posterior standard deviation for a single location parameter when analysed with a $t$-distribution, compared to when analysed with a normal distribution. Each point represents a dataset with $N = 100$ datapoints. The black dashed line shows the expected ratio of posterior standard deviations as calculated using the Fisher information of the $t$-likelihood.

## 3.1 Simulation-based calibration

Simulation-based calibration tests (Cook et al. 2006; Talts et al. 2018) are used to diagnose sampling issues. Parameters of interest $\theta_0$ are drawn from the model prior $\pi(\theta)$, and a dataset is drawn following the prescription of the model. We can then run a single MCMC chain on this dataset until we have drawn $L$ independent samples $\{\tilde{\theta}_i\}$ (see Talts et al. (2018) for a discussion of independence criteria). We then compute the rank statistic

$$r(\theta_0, \tilde{\boldsymbol{\theta}}_{\text{MCMC}}) = \sum_{i=1}^{L} \mathbb{I}(\tilde{\theta}_i < \theta_0), \tag{22}$$

where $\mathbb{I}(\cdot)$ is the indicator function which evaluates to 1 if its (logical) argument is true and 0 if it is false. If the generative model matches the sampling model and the sampling algorithm correctly draws from the posterior, the rank statistic will be uniformly distributed across the integers $[0, L]$. We can, therefore, repeat this process many times to build a histogram of rank statistics; any shortfalls in the sampling algorithm will manifest as deviations from uniformity.

While this procedure has been developed to diagnose sampling issues, we can also use the method to build a heuristic measure of how robust a model is to mis-specification. If the generative model and the sampling model no longer match, we would expect deviations from uniformity in the rank statistic histogram. We can compare how well different sampling models deal with model mis-specifcation by building the rank statistic histogram and assessing the deviation from uniformity.

As our priors are defined in a scaled space, our simulation-based calibration tests are also conducted in this scaled space; in this way, we are solely testing the performance of the MCMC models, and not of the scaling before fitting the models.

## 3.2 Fixed-value calibration

For our fixed-value calibration tests, we simulate datasets from known models with true fixed values for our regression parameters. The purpose of these tests is to demonstrate accurate recovery of these values, and to compare with the results given when using normal distributions.

We generate multiple datasets with the same ground-truth parameters to verify the results across multiple runs, producing plots of the composite cumulative distribution function across each dataset. In each instance, full specifications of the data models are given in Appendix C.

We aim to do a full end-to-end test of our inference pipeline in the fixed-value calibration tests; therefore, these datasets are generated in the unscaled space and are a test not only of the sampler but also of the scaling.

## 4 RESULTS ON SIMULATED DATASETS

In the previous section, we proposed a general-purpose, robust statistical model for linear regression; in this section, we investigate the performance of the model on a series of simulated datasets with known parameters. Code that reproduces the datasets in this section is available online[2].

### 4.1 $t$-distributed data

For our first test, we start with a dataset that matches our model perfectly (i.e. $\mathcal{P}_{\text{int}} = t_\nu, \mathcal{P}_{\text{obs}} = \mathcal{N}$), and choose $K = 1$ independent variables. The simulation-based calibration tests indicate that the rank statistic is consistent with being uniformly distributed (see Figure 5), and, therefore, that our inference procedure is working as expected.

For our fixed-value tests, we drew $N = 20$ datapoints with $(\alpha, \beta, \sigma_{68}, \nu) = (3, 2, 0.1, 3)$; the full model used to generate the data is specified in Appendix C1. The chosen value of $\nu = 3$ corresponds to an outlier fraction of $\omega(\nu = 3) \approx 5.8\%$. As shown in Figure 6, we recover the values of the parameters that were used to generate the dataset; these values are consistent with a model where the shape parameter is fixed to $\nu = 3$.

We then generated 400 datasets from the same ground-truth parameters, and ran MCMC against each dataset. 95% highest posterior density credible intervals were constructed for each run; the true parameter values were contained in these credible intervals 97%, 96% and 98% of the time for the intercept $\alpha$, the slope $\beta$ and the intrinsic scatter $\sigma_{68}$. For the nuisance parameter $\nu$, the credible intervals contained the true parameter value across all runs; this may be caused by this parameter being only weakly constrained for these datasets.

### 4.2 Normally-distributed data

In this test, we compare $t$-cup with an equivalent model that employs normal distributions to check that:

(i) $t$-cup reduces to a normal model in the absence of outliers

(ii) $t$-cup gives less biased results when an extreme outlier is present.

We conducted a simulation-based calibration test, where datasets were generated using normal distributions throughout (i.e. $\mathcal{P}_{\text{int}} = \mathcal{P}_{\text{obs}} = \mathcal{N}$) and with a single independent variable (i.e. $K = 1$). The results (shown in Figure 7) indicate that, while the estimates of the intercept and slope are significantly biased when analysed by the normal model (as indicated by the rank statistic's distribution
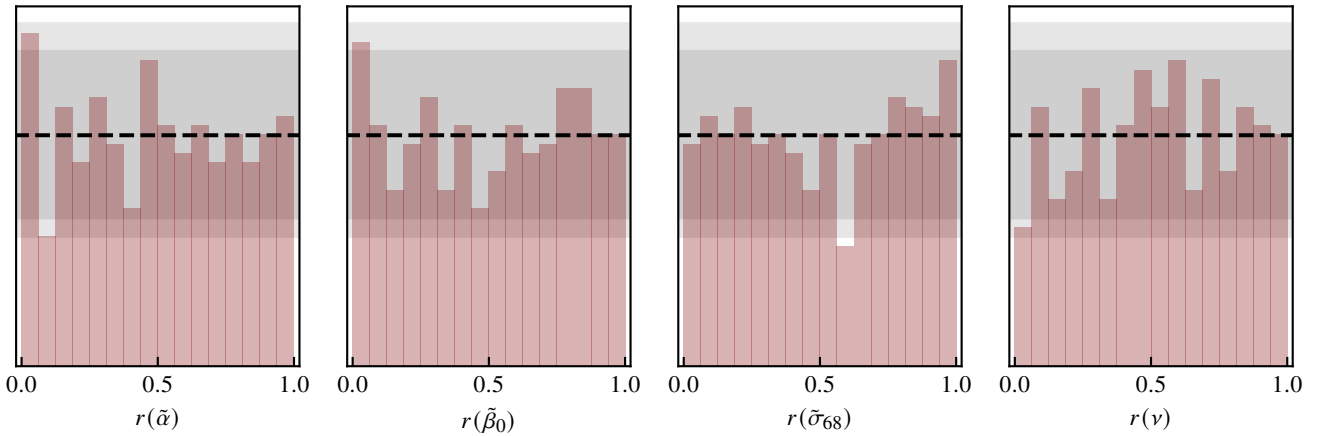
---

[2] https://github.com/wm1995/tcup-paper

**Figure 5.** SBC runs for the t distribution under *t*-cup. If the inference procedure is working as expected, the histograms for each parameter should be distributed uniformly (as indicated by the black dashed line). The dark (light) grey regions correspond to the 94% (98%) confidence interval of uniformity (i.e. we expect one histogram bin per panel (figure) to lie outside of this range).
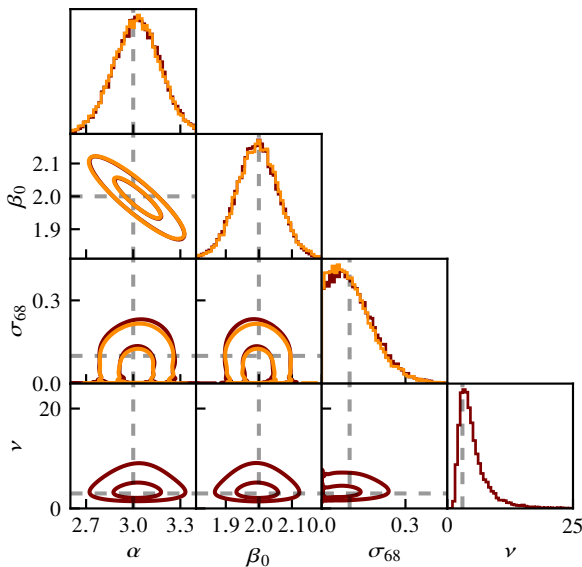


**Figure 6.** The posterior for each of the regression coefficients and the shape parameter $\nu$ for the *t*-distributed dataset, with ground-truth values indicated by the black dashed lines. Constraints from *t*-cup (red) are compared with those derived with $\nu$ fixed to the true value (orange). Contours indicate 39.3% and 86.5% highest posterior density regions, corresponding to $1\sigma$ and $2\sigma$ contours for a bivariate normal distribution.

deviating from normal), the estimates derived by *t*-cup are unbiased. Additionally, while both *n*-cup and *t*-cup systematically overestimate the intrinsic scatter $\sigma_{68}$, the bias is much smaller for *t*-cup than for *n*-cup.

A dataset of $N = 12$ points was generated with $(\alpha, \beta, \sigma_{int}) = (3, 2, 0.2)$, and one of the points was modified to be a $\sim 20\,\sigma$ outlier (the full generative model is given in Appendix C2). While such an extreme outlier could easily be identified and removed, the purpose

here is to demonstrate that *t*-cup does not require this to obtain sensible inferences.

Figure 8 illustrates how the estimates for the true parameter models are biased in the normal model, but less affected in the *t*-cup model.

In Figure 9, we compare the constraints on parameters derived under the normal model (including and excluding the outlier from the dataset), and the *t*-cup model (including the outlier only). The constraints from the *t*-cup model including the outlier are consistent with those derived under the normal model when the outlier is excluded, obviating the need to remove the outlier manually. While this outlier is particularly extreme, this example illustrates the utility of the *t*-cup model in datasets with outliers.

For completeness, in Figure 10 we illustrate that the *t*-cup model recovers consistent constraints regardless of whether the outlier is included or excluded.

We then generated 400 datasets using the same procedure, and combined posterior samples from each run to build an effective cumulative distribution function across all runs for both the normal model and for *t*-cup. The results (illustrated in Figure 11) indicate that constraints under *t*-cup are significantly less biased than those calculated under the normal model.

### 4.3 Two-dimensional normal mixture model with outliers

This test is designed to further explore how *t*-cup performs when the model is misspecified. In this case, we are looking at a normally-distributed population which has a 10% contamination rate with another normally-distributed population of the same mean but 10 times the standard deviation. Equivalently, this test can be thought of as investigating how the model performs when there is a significant fraction of outliers.

The intrinsic scatter distribution $\mathcal{P}_{int}$ is a mixture of two normal distributions with zero mean; 90% of points are drawn from a core distribution with standard deviation $\sigma_{int}$, and 10% of points are drawn from an outlier distribution with standard deviation $10\sigma_{int}$. The observation distribution $\mathcal{P}_{obs}$ is a normal distribution. For fixed-value tests, the true values of the regression parameters were fixed to $(\alpha, \beta_0, \beta_1, \sigma_{int}) = (2, 3, 1, 0.4)$. The full generative model for the fixed-value tests is given in Appendix C3.

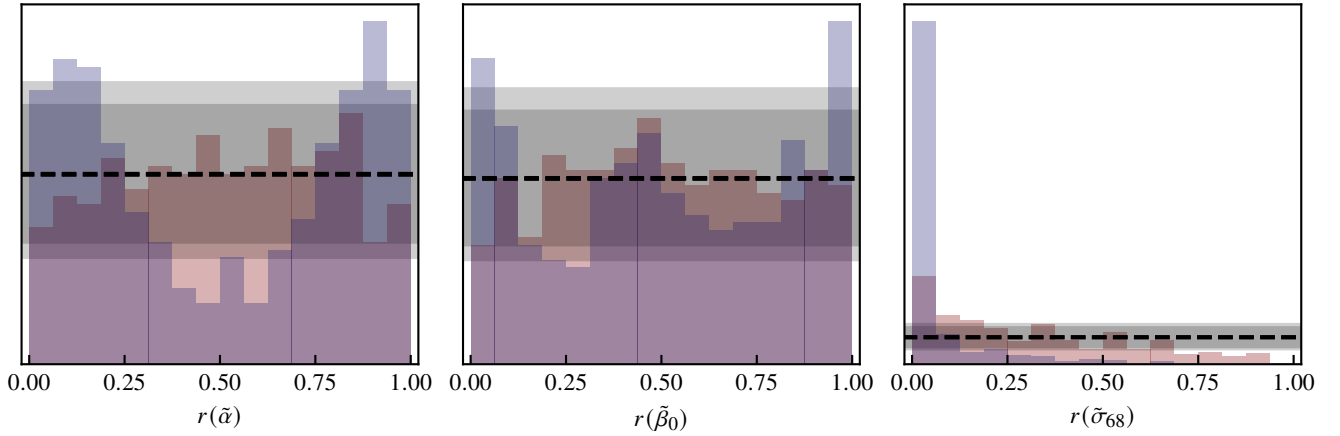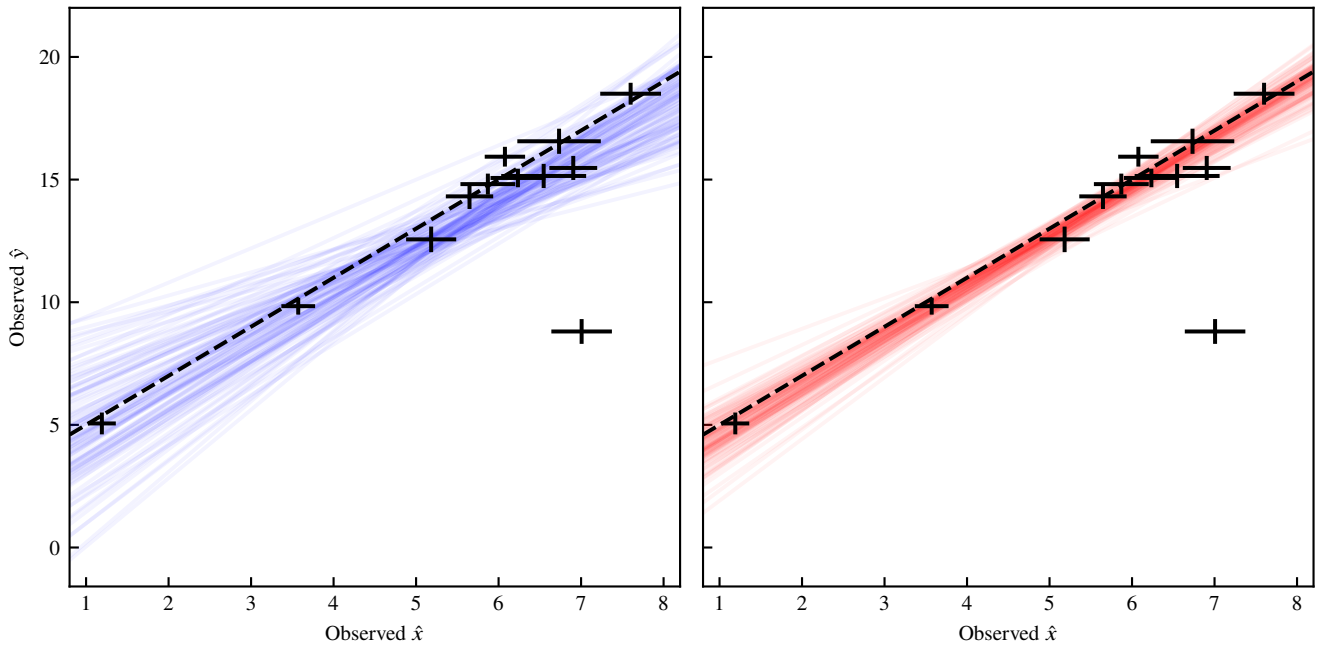**Figure 7.** A rank statistic histogram for 400 simulation-based calibration runs for the outlier dataset under $n$-cup (blue) and $t$-cup (red). If there were a perfect match between the generative model and the inference model, the histograms for each parameter will be distributed uniformly (as indicated by the black dashed line). The dark (light) grey regions correspond to the 94% (98%) confidence interval of uniformity (i.e. we expect one histogram bin per panel (figure) to lie outside of this range).



**Figure 8.** 100 draws from the posterior of regression lines from the normal model (left panel, blue) and $t$-cup (right panel, red). The dataset is illustrated by the black points, with the ground-truth regression line illustrated by the black dashed line.

The results (see Figure 12) show that, while the constraints from the normal model are significantly biased for this generative distribution, $t$-cup is able to recover the correct parameter values.

### 4.4 Laplace-distributed data

As the underlying sampling distributions can never be known with certainty it is important to test the method on data generated on multiple distributions that do not fall into the family of $t$-distributions. Here we use a Laplace distribution as $\mathcal{P}_{\text{int}}$, giving another example

of performance under explicit model mis-specification. The Laplace distribution has probability density

$$\text{Laplace}\,(x; \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right), \tag{23}$$

with scale parameter $b$ related to $\sigma_{68}$ as $b(\sigma_{68}) = -\sigma_{68}/\ln\left(1 - \text{erf}(1/\sqrt{2})\right) \approx 0.87\,\sigma_{68}$.

We generate $N = 25$ datapoints from a model with Laplacian intrinsic scatter (i.e. $\mathcal{P}_{\text{int}} = \text{Laplace}$), normally distributed observation noise (i.e. $\mathcal{P}_{\text{obs}} = \mathcal{N}$), and $K = 1$ independent variables. For

**Figure 9.** The posterior for each of the regression coefficients under the normal model with the outlier included (dark blue, solid) and excluded (light blue, dashed), and for $t$-cup with the outlier included (dark red, solid). Ground-truth values are indicated by the black dashed lines.



**Figure 10.** The posterior for each of the regression coefficients and the shape parameter $\nu$ under the $t$-cup model with the outlier included (dark red, solid) and excluded (light red, dashed). Ground-truth values are indicated by the black dashed lines.

**Table 1.** Estimates for the intercept, slope, and intrinsic scatter inferred for the relationship between Eddington ratio $L/L_{\rm Edd}$ and X-ray spectral index $\Gamma$ for the sample of quasars analysed by Kelly (2007). The parameter estimates reported in Kelly (2007) are the posterior median and "a robust estimate of the standard deviation"; for linmix and $t$-cup, we report the posterior median and an estimate of the standard deviation as $\sigma = 1.4826$ MAD, where MAD is the median absolute deviation.

|  | Kelly (2007) | linmix | $t$-cup |
|---|---|---|---|
| Intercept $\alpha$ | $3.12 \pm 0.41$ | $3.18 \pm 0.48$ | $3.62 \pm 0.26$ |
| Slope $\beta$ | $1.35 \pm 0.54$ | $1.40 \pm 0.60$ | $2.00 \pm 0.33$ |
| Int. scatter $\sigma_{68}$ | $0.26 \pm 0.11$ | $0.25 \pm 0.12$ | $0.08 \pm 0.07$ |
| Outlier fraction $\omega$ | — | — | $0.04 \pm 0.03$ |

the fixed-value tests, we set $(\alpha, \beta, \sigma_{\rm int}) = (-1, 0.8, 0.2)$. The full generative model can be found in Appendix C4.

As we can see in Figure 13, while both $n$-cup and $t$-cup successfully recover the gradient and intercept that were used to generate the dataset, $n$-cup overpredicts the true intrinsic scatter when $t$-cup constrains it accurately. To confirm this, we generated 400 datasets using the same procedure, and analysed each with both $n$-cup and $t$-cup. The results (illustrated in Figure 14) show the same pattern – that $t$-cup is able to accurately constrain the intrinsic scatter, while the estimate from the normal model is biased high.

## 5 DEMONSTRATION ON REAL DATA

Having seen $t$-cup's performance on simulated data in the previous section, we now compare the performance of the $t$-cup model with a generic astronomical Bayesian linear regression model (LINMIX_ERR; Kelly 2007) and a tailored approach using $t$-distributions (Park et al. 2017).

### 5.1 LINMIX_ERR

We use the same dataset that is used in Section 8 of Kelly (2007) – a dataset of $N = 39$ quasars with measured Eddington ratio $L/L_{\rm Edd}$ and X-ray spectral index $\Gamma$. Performance is compared with a Python implementation[3] of the original LINMIX_ERR paper proposed by Kelly – see Figure 15.

While the parameter estimates are broadly consistent (see Table 1), the posterior distributions in Figure 16 can be seen to differ, with the posterior inferred by LINMIX being more diffuse than that inferred by $t$-cup. The tighter constraints of $t$-cup suggest that the inference presented by Kelly (2007) may be biased by the enforced assumption of normally-distributed data, which may not be accurate. The explicit assumption of normality in LINMIX is in tension with the outlier fraction estimated by $t$-cup (68% CI $\omega = 0.05^{+0.04}_{-0.01}$). This showcases that analysing data with robust procedures gives materially different answers on real-world data, and highlights the importance of carefully considering the implications when assuming normality.

### 5.2 Park et al. (2017)

Park et al. (2017) presents a bespoke $t$-distribution-based linear regression model for estimating SMBH mass; this allows us to compare the results of their bespoke model with our generic one. Their dataset consists of $N = 31$ AGN with reverberation-mapped mass estimates, $M_{\rm BH}$, C IV emission line width, $\Delta V$, and continuum luminosities,
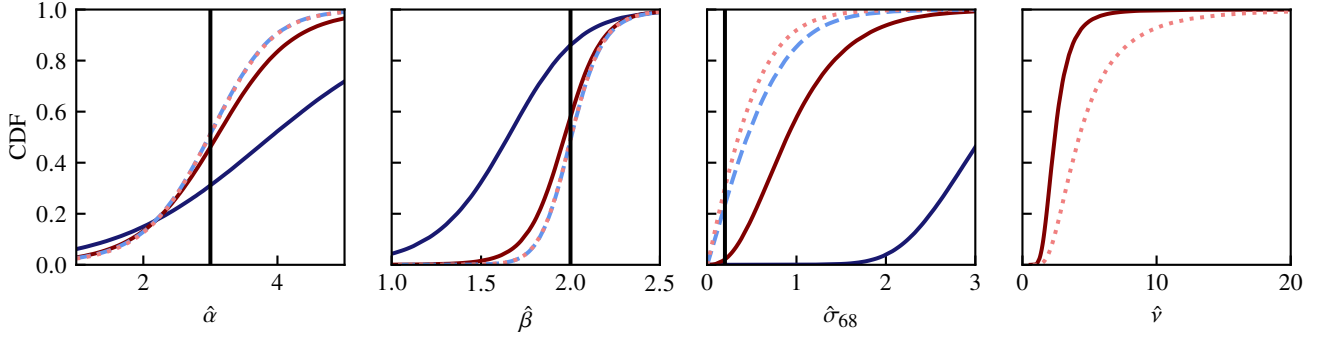
---

[3] https://github.com/jmeyers314/linmix

**Figure 11.** The combined cumulative distribution function of the regression parameters for 400 normal datasets with an outlier under the normal model (dark blue, solid) and *t*-cup model (dark red, solid), and with the outlier removed for the normal model (light blue, dashed) and *t*-cup model (light red, dashed).
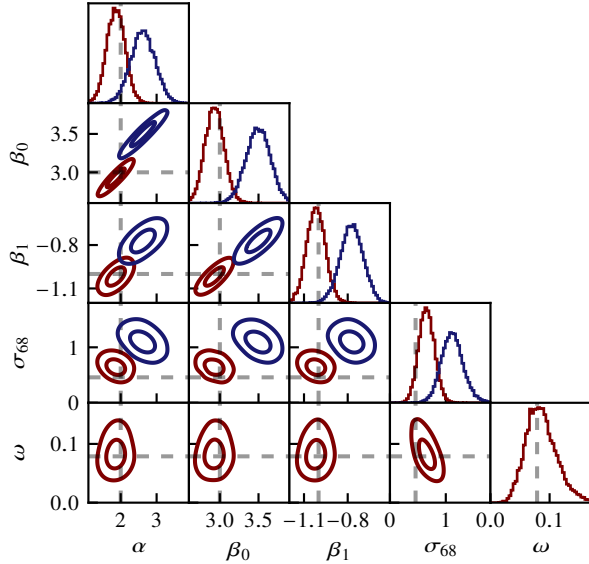


**Figure 12.** The posterior for each of the regression coefficients and the outlier fraction $\omega$ for the normal mixture model dataset for *n*-cup (blue) and *t*-cup (red), with ground-truth values indicated by the grey dashed lines.



**Figure 13.** The posterior for each of the regression coefficients and the shape parameter $\nu$ for the Laplace-distributed dataset under the normal model (blue) and *t*-cup (red), with ground-truth values indicated by the black dashed lines.

$\lambda L_\lambda$. Three different C IV line width measurements are compared in Park et al. (2017): the full-width at half maximum (FWHM), the line dispersion, $\sigma_l$, and the median absolute deviation. This dataset is used to fit the regression relation

$$\log_{10}\left(\frac{M_{\mathrm{BH}}}{10^8 M_\odot}\right) = \alpha + \beta \log_{10}\left(\frac{\lambda L_\lambda}{10^{44}\mathrm{erg\ s^{-1}}}\right) + \gamma \log_{10}\left(\frac{\Delta V}{10^3 \mathrm{km\ s^{-1}}}\right). \quad (24)$$

The *t*-cup posterior for parameters $\alpha$, $\beta$ and $\gamma$, as well as intrinsic scatter $\sigma_{68}$ and shape parameter $\nu$, is shown in Figure 17.

The constraints on the intercept, $\alpha$, and slopes, $\beta$ and $\gamma$, are consistent with those derived by Park et al. (2017) for all three measures of emission line width — see Table 2. Park et al. (2017) predicts a systematically higher intrinsic scatter than that predicted by *t*-cup, however, these figures cannot be compared directly, as this may be influenced by different prior choice for $\nu$.

## 6 CONCLUSIONS

We have presented a general-purpose approach to linear regression, implemented as *t*-cup, that is robust to model mis-specification, with a model laid out in Section 2. In Section 4, we demonstrated that the model recovers constraints consistent with those used to generate the datasets, including in cases where there is a mismatch between the generative model used to create the dataset and our regression model. In Section 5, we compared the method to Kelly (2007); Park et al. (2017) on real-world data, illustrating that the models derive consistent constraints.

It may be fruitful to re-examine some of the priors assumed in the *t*-cup model; while we focused on producing a model that was applicable to a large range of datasets, *t*-cup predicted lower intrinsic scatter than both Kelly (2007) and Park et al. (2017). While results between the three are not directly comparable as a result of the
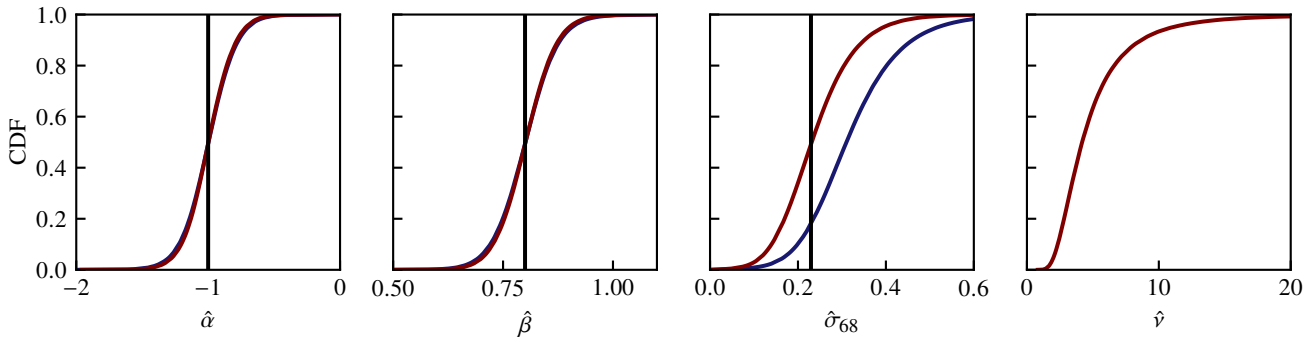
**Figure 14.** The combined cumulative distribution function of the regression parameters for 400 Laplace datasets under the normal model (blue) and $t$-cup model (red).
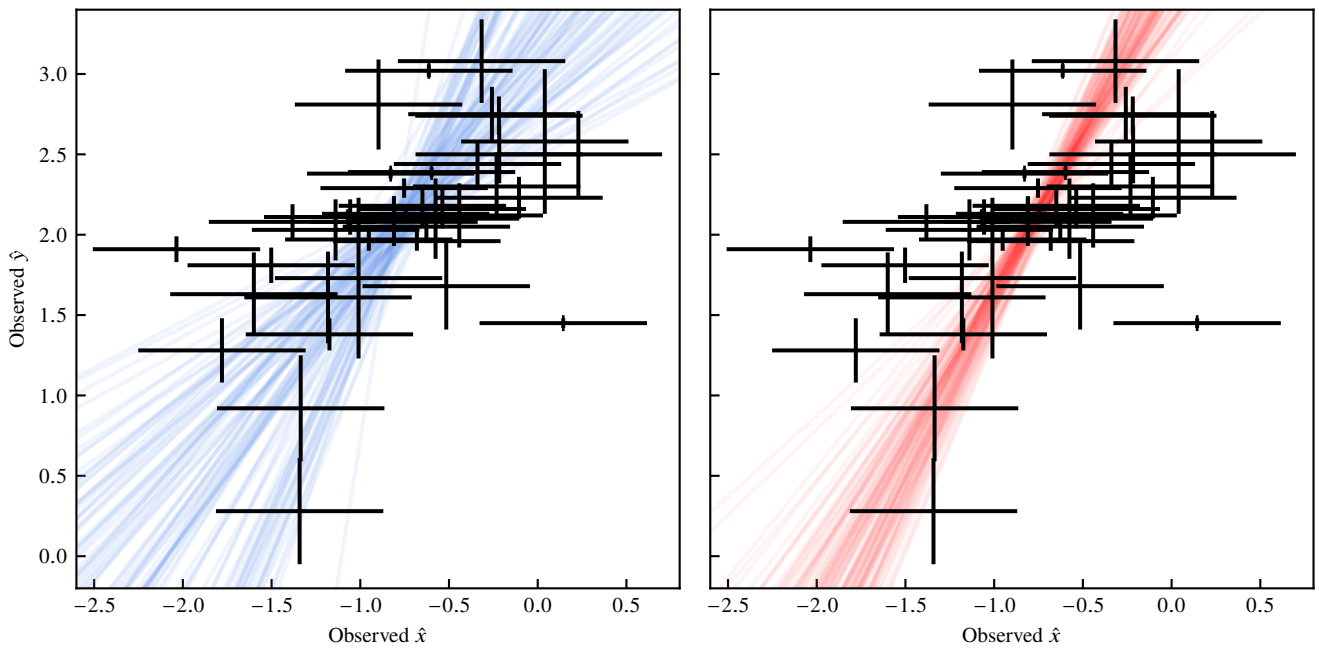


**Figure 15.** 100 draws from the posterior of regression lines from LINMIX_ERR (left panel, blue) and $t$-cup (right panel, red) for the Kelly (2007) dataset.

**Table 2.** A comparison between the constraints on the parameters in Equation 24 derived by Park et al. (2017) and using $t$-cup.

|  | FWHM | | $\sigma_l$ | | MAD | |
|---|---|---|---|---|---|---|
|  | Park et al. (2017) | $t$-cup | Park et al. (2017) | $t$-cup | Park et al. (2017) | $t$-cup |
| Intercept $\alpha$ | $7.54^{+0.26}_{-0.27}$ | $7.51^{+0.22}_{-0.22}$ | $6.90^{+0.35}_{-0.34}$ | $6.91^{+0.30}_{-0.31}$ | $7.15^{+0.24}_{-0.25}$ | $7.15^{+0.22}_{-0.22}$ |
| Slope $\beta$ | $0.45^{+0.08}_{-0.08}$ | $0.43^{+0.06}_{-0.06}$ | $0.44^{+0.07}_{-0.07}$ | $0.43^{+0.05}_{-0.06}$ | $0.42^{+0.07}_{-0.06}$ | $0.42^{+0.05}_{-0.06}$ |
| Slope $\gamma$ | $0.50^{+0.55}_{-0.53}$ | $0.58^{+0.44}_{-0.45}$ | $1.66^{+0.65}_{-0.66}$ | $1.66^{+0.57}_{-0.58}$ | $1.65^{+0.61}_{-0.62}$ | $1.65^{+0.55}_{-0.55}$ |
| Int. scatter $\sigma$ | $0.16^{+0.10}_{-0.08}$ | $0.10^{+0.03}_{-0.10}$ | $0.12^{+0.09}_{-0.06}$ | $0.08^{+0.02}_{-0.08}$ | $0.12^{+0.09}_{-0.06}$ | $0.07^{+0.02}_{-0.07}$ |

different assumptions in each model, this difference could suggest that the prior on $\sigma_{68}$ has too much density near $\sigma_{68} = 0$. In addition, setting the prior on $\{\mathbf{x}_i\}$ using extreme deconvolution, while effective, is not theoretically motivated; a prior similar to that presented in Bartlett & Desmond (2023) may be more appropriate in this case.

In our next paper, we will apply the robust inference techniques presented here to exploring the single-epoch mass estimators used in estimating the masses of high-redshift quasars.
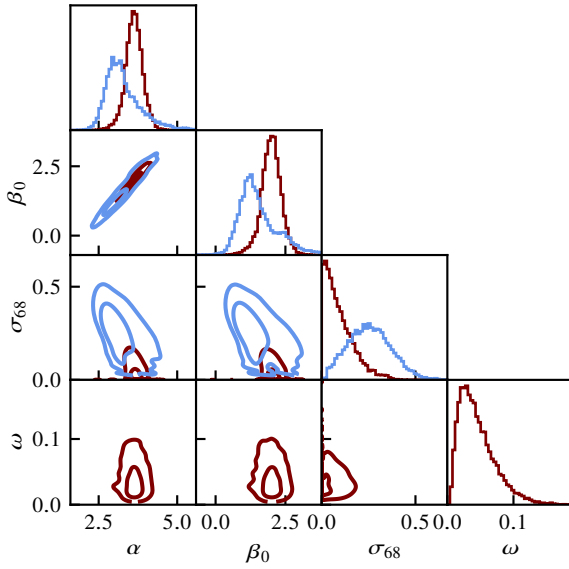
**Figure 16.** The posterior distributions for the data from Kelly (2007) under the linmix (blue) and $t$-cup (red) models.
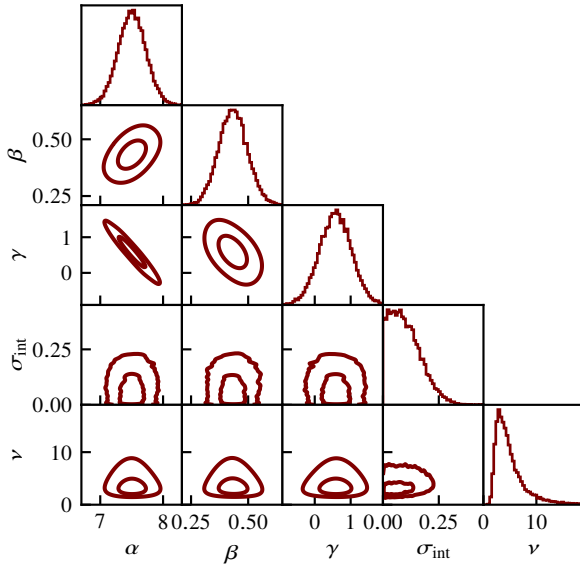


**Figure 17.** The posterior distributions for the data from Park et al. (2017) under the $t$-cup model. This figure is directly comparable with Figure 9 from Park et al. (2017).

## DATA AVAILABILITY

The data in this paper is sourced from Kelly (2007); Park et al. (2017). Scripts to generate the simulated datasets used to validate the model in Section 4 are available from the GitHub repository for this paper, at https://github.com/wm1995/tcup-paper.

## REFERENCES

Aitkin, M. & Tunnicliffe Wilson, G., 1980. Mixture models, outliers, and the em algorithm, *Technometrics*, **22**, 325–331.

Akritas, M. G. & Bershady, M. A., 1996. Linear Regression for Astronomical Data with Measurement Errors and Intrinsic Scatter, *ApJ*, **470**, 706.

Andreon, S., 2020. Evidence for radially independent size growth of early-type galaxies in clusters, *A&A*, **640**, A34.

Andreon, S. & Hurn, M., 2013. Measurement errors and scaling relations in astrophysics: a review, *Statistical Analysis and Data Mining: The ASA Data Science Journal*, **9**(1), 15–33.

Andreon, S. & Weaver, B., 2015. *Bayesian Methods for the Physical Sciences. Learning from Examples in Astronomy and Physics.*, vol. 4, Springer Cham.

Andreon, S., Puddu, E., de Propris, R., & Cuillandre, J. C., 2008. Galaxy evolution in the high-redshift, colour-selected cluster RzCS 052 at z = 1.02, *MNRAS*, **385**(2), 979–985.

Andrews, D. F. & Mallows, C. L., 1974. Scale Mixtures of Normal Distributions, *Journal of the Royal Statistical Society: Series B (Methodological)*, **36**(1), 99–102.

Andrews, S. M., Rosenfeld, K. A., Kraus, A. L., & Wilner, D. J., 2013. The Mass Dependence between Protoplanetary Disks and their Stellar Hosts, *ApJ*, **771**(2), 129.

Avelino, A., Friedman, A. S., Mandel, K. S., Jones, D. O., Challis, P. J., & Kirshner, R. P., 2019. Type Ia Supernovae Are Excellent Standard Candles in the Near-infrared, *ApJ*, **887**(1), 106.

Bailey, D. C., 2017. Not normal: the uncertainties of scientific measurements, *Royal Society Open Science*, **4**(1), 160600.

Bartlett, D. J. & Desmond, H., 2023. Marginalised Normal Regression: Unbiased curve fitting in the presence of x-errors, *The Open Journal of Astrophysics*, **6**, 42.

Bentz, M. C., Denney, K. D., Grier, C. J., Barth, A. J., Peterson, B. M., Vestergaard, M., Bennert, V. N., Canalizo, G., De Rosa, G., Filippenko, A. V., Gates, E. L., Greene, J. E., Li, W., Malkan, M. A., Pogge, R. W., Stern, D., Treu, T., & Woo, J.-H., 2013. The Low-luminosity End of the Radius-Luminosity Relationship for Active Galactic Nuclei, *ApJ*, **767**(2), 149.

Berger, J. O., Moreno, E., Pericchi, L. R., Bayarri, M. J., Bernardo, J. M., Cano, J. A., la Horra, J. D., Martín, J., Ríos-Insúa, D., Betrò, B., Dasgupta, A., Gustafson, P., Wasserman, L., Kadane, J. B., Srinivasan, C., Lavine, M., O'Hagan, A., Polasek, W., Robert, C. P., Goutis, C., Ruggeri, F., Salinetti, G., & Sivaganesan, S., 1994. An overview of robust bayesian analysis, *Test*, **3**(1), 5–124.

Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P. A., Horsfall, P., & Goodman, N. D., 2019. Pyro: Deep universal probabilistic programming, *J. Mach. Learn. Res.*, **20**, 28:1–28:6.

Bovy, J., Hogg, D. W., & Roweis, S. T., 2011. Extreme deconvolution: Inferring complete distribution functions from noisy, heterogeneous and incomplete observations, *Annals of Applied Statistics*, **5**(2), 1657–1677.

Box, G. E. P. & Tiao, G. C., 1968. A bayesian approach to some outlier problems, *Biometrika*, **55**(1), 119–129.

Chung, Y., Rabe-Hesketh, S., Dorie, V., Gelman, A., & Liu, J., 2013. A nondegenerate penalized likelihood estimator for variance parameters in multilevel models, *Psychometrika*, **78**(4), 685–709.

Cook, S. R., Gelman, A., & Rubin, D. B., 2006. Validation of software for bayesian models using posterior quantiles, *Journal of Computational and Graphical Statistics*, **15**(3), 675–692.

Cox, R. T., 1946. Probability, Frequency and Reasonable Expectation, *American Journal of Physics*, **14**(1), 1–13.

Cramér, H., 1946. *Mathematical methods of statistics*, Princeton mathematical series ; 9, Princeton University Press, Princeton.

Ding, P., 2014. Bayesian robust inference of sample selection using selection-t models, *Journal of Multivariate Analysis*, **124**, 451–464.

Efstathiou, G., 2014. H₀ revisited, *MNRAS*, **440**(2), 1138–1152.

Feeney, S. M., Mortlock, D. J., & Dalmasso, N., 2018. Clarifying the Hubble constant tension with a Bayesian hierarchical model of the local distance ladder, *MNRAS*, **476**(3), 3861–3882.

Feigelson, E. D., de Souza, R. S., Ishida, E. E. O., & Jogesh Babu, G., 2021. 21st Century Statistical and Computational Challenges in Astrophysics, *Annual Review of Statistics and Its Application*, **8**, 493–517.

Ferrarese, L. & Merritt, D., 2000. A Fundamental Relation between Super-massive Black Holes and Their Host Galaxies, *ApJ*, **539**(1), L9–L12.

Gebhardt, K., Bender, R., Bower, G., Dressler, A., Faber, S. M., Filippenko, A. V., Green, R., Grillmair, C., Ho, L. C., Kormendy, J., Lauer, T. R., Magorrian, J., Pinkney, J., Richstone, D., & Tremaine, S., 2000. A Relationship between Nuclear Black Hole Mass and Galaxy Velocity Dispersion, *ApJ*, **539**(1), L13–L16.

Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B., 2013. *Bayesian Data Analysis*, Chapman and Hall/CRC.

Ghosh, J. K., Delampady, M., & Samanta, T., 2006. *An Introduction to Bayesian Analysis*, Springer-Verlag.

Hubble, E., 1929. A Relation between Distance and Radial Velocity among Extra-Galactic Nebulae, *Proceedings of the National Academy of Science*, **15**(3), 168–173.

Jing, T. & Li, C., 2024. Regression for Astronomical Data with Realistic Distributions, Errors and Non-linearity, *arXiv e-prints*, p. arXiv:2411.08747.

Jontof-Hutter, D., Ford, E. B., Rowe, J. F., Lissauer, J. J., Fabrycky, D. C., Van Laerhoven, C., Agol, E., Deck, K. M., Holczer, T., & Mazeh, T., 2016. Secure Mass Measurements from Transit Timing: 10 Kepler Exoplanets between 3 and 8 $M_\oplus$ with Diverse Densities and Incident Fluxes, *ApJ*, **820**(1), 39.

Juárez, M. A. & Steel, M. F. J., 2010. Non-gaussian dynamic bayesian modelling for panel data, *Journal of Applied Econometrics*, **25**(7), 1128–1154.

Kelly, B. C., 2007. Some Aspects of Measurement Error in Linear Regression of Astronomical Data, *ApJ*, **665**(2), 1489–1506.

Kelly, B. C., Vestergaard, M., & Fan, X., 2009. Determining Quasar Black Hole Mass Functions from their Broad Emission Lines: Application to the Bright Quasar Survey, *ApJ*, **692**(2), 1388–1410.

Knuth, K. H. & Skilling, J., 2010. Foundations of Inference, *arXiv e-prints*, p. arXiv:1008.4831.

Leavitt, H. S. & Pickering, E. C., 1912. Periods of 25 Variable Stars in the Small Magellanic Cloud., *Harvard College Observatory Circular*, **173**, 1–3.

Leistedt, B., Mortlock, D. J., & Peiris, H. V., 2016. Hierarchical Bayesian inference of galaxy redshift distributions from photometric surveys, *MNRAS*, **460**(4), 4258–4267.

Mantz, A. B., 2016. A Gibbs sampler for multivariate linear regression, *MNRAS*, **457**(2), 1279–1288.

McConnell, N. J. & Ma, C.-P., 2013. Revisiting the Scaling Relations of Black Hole Masses and Host Galaxy Properties, *ApJ*, **764**(2), 184.

Neal, R. M., 2003. Slice sampling, *The Annals of Statistics*, **31**(3), 705 – 767.

Park, D., Barth, A. J., Woo, J.-H., Malkan, M. A., Treu, T., Bennert, V. N., Assef, R. J., & Pancoast, A., 2017. Extending the Calibration of C IV-based Single-epoch Black Hole Mass Estimators for Active Galactic Nuclei, *ApJ*, **839**(2), 93.

Phan, D., Pradhan, N., & Jankowiak, M., 2019. Composable effects for flexible and accelerated probabilistic programming in numpyro, *arXiv preprint arXiv:1912.11554*.

Press, W. H., Teukolsky, S. A., Vetterling, W. T., & Flannery, B. P., 1992. *Numerical recipes in C. The art of scientific computing*, Cambridge University Press.

Rao, C. R., 1945. Information and the accuracy attainable in the estimation of statistical parameters, *Bull. Calcutta Math. Soc.*, **37**, 81–91.

Riess, A. G., Macri, L., Casertano, S., Lampeitl, H., Ferguson, H. C., Filippenko, A. V., Jha, S. W., Li, W., & Chornock, R., 2011. A 3% Solution: Determination of the Hubble Constant with the Hubble Space Telescope and Wide Field Camera 3, *ApJ*, **730**(2), 119.

Sereno, M., 2016. A Bayesian approach to linear regression in astronomy, *MNRAS*, **455**(2), 2149–2162.

Sestovic, M., Demory, B.-O., & Queloz, D., 2018. Investigating hot-Jupiter inflated radii with hierarchical Bayesian modelling, *A&A*, **616**, A76.

Sivia, D. S. & Skilling, J., 2006. *Data analysis : a Bayesian tutorial*, Oxford science publications, Oxford University Press, Oxford ;, 2nd edn.

Tak, H., Ellis, J. A., & Ghosh, S. K., 2019. Robust and accurate inference via a mixture of gaussian and student's t errors, *Journal of Computational and Graphical Statistics*, **28**(2), 415–426.

Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A., 2018. Validating Bayesian Inference Algorithms with Simulation-Based Calibration, *arXiv e-prints*, p. arXiv:1804.06788.

Tremaine, S., Gebhardt, K., Bender, R., Bower, G., Dressler, A., Faber, S. M., Filippenko, A. V., Green, R., Grillmair, C., Ho, L. C., Kormendy, J., Lauer, T. R., Magorrian, J., Pinkney, J., & Richstone, D., 2002. The Slope of the Black Hole Mass versus Velocity Dispersion Correlation, *ApJ*, **574**(2), 740–753.

Van Horn, K. S., 2003. Constructing a logic of plausible inference: a guide to cox's theorem, *International Journal of Approximate Reasoning*, **34**(1), 3–24.

## APPENDIX A: PRIOR CHOICE FOR SHAPE PARAMETER

As our model relies on Student's *t*-distributions, we review notation and priors used by previous works and justify our reasoning for our prior choice. One approach is to adopt a fixed value of $\nu$, which is equivalent to setting a Dirac delta function prior: a common choice is $\nu = 4$ (e.g. Berger et al. 1994; Gelman et al. 2013). Another approach is to adopt a more flexible approach by allowing $\nu$ to vary (e.g. Juárez & Steel 2010; Gelman et al. 2013; Ding 2014; Park et al. 2017; Feeney et al. 2018).

We sought a flexible prior for this work that could reduce to a (nearly) normal distribution, but had sufficient flexibility to accommodate heavy-tailed distributions as well. Priors meeting this criterion include:

- priors of the form $\nu \sim \Gamma(\alpha, \beta)$ for some shape parameter, $\alpha$, and rate parameter, $\beta$ – e.g. Juárez & Steel (2010) uses $\{\alpha = 2, \beta = 0.1\}$; Ding (2014) uses $\{\alpha = 1, \beta = 0.1\}$
- A uniform prior in $\frac{1}{\nu} \sim U(0, 1)$ (Gelman et al. 2013)
- A uniform prior[4] in distribution peak height relative to a normal distribution, $t$, such that

$$t \equiv \sqrt{\frac{2}{\nu}} \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \sim U(0, 1). \tag{A1}$$

These priors, along with the prior we adopted, are illustrated in Figure A1. The cumulative distribution functions for the priors in terms of outlier fraction, $\omega$, are illustrated in Figure A2.

When testing different priors, we found that priors with significant density in the range $0 < \nu < 1$ led to sampling difficulties as $\nu$ approached 0; these low values of $\nu \sim 0.1$ can correspond to outliers that are more than 20 orders of magnitude larger than $\sigma$. In theory, these unphysical regions of parameter space ought to be excluded during the process of inference. However, the use of HMC in such cases can lead to divergences in the sampling process or inefficient sampling as the sampler struggles with regions of high curvature — see the discussion of this phenomenon in Neal (2003).

Another option would be to place a prior is on outlier fraction $\omega$. It could be argued that the term outlier loses its meaning when the

---

[4] Strictly, Feeney et al. (2018) approximate this prior with the closed form prior on shape parameter, $\nu$, of

$$\mathcal{P}(\nu) \propto \frac{\Theta(\nu)}{\left(\left(\frac{\nu}{\nu_0}\right)^{1/(2a)} + \left(\frac{\nu}{\nu_0}\right)^{2/a}\right)^a}$$

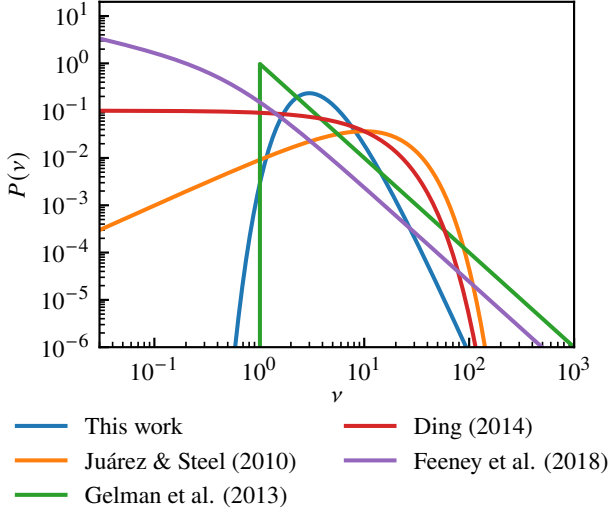where $\Theta(\cdot)$ is the Heaviside step function, and $\nu_0$ and $a$ are constants with values $\sim 0.55$ and $\sim 1.2$ respectively.

**Figure A1.** A selection of different priors on $\nu$ in works that have used $t$-distributions.



**Figure A2.** The CDF of priors on $\nu$ used in works that have used $t$-distributions. The priors are expressed in terms of outlier fraction, $\omega$, as defined in Equation 12.

majority of a dataset is composed of so-called "outliers"; therefore, a natural choice of prior might be a uniform distribution ranging from the normally-distributed outlier fraction of ~0.00270 to this "maximum" outlier fraction of 0.5, corresponding to $\nu \sim 0.302$. This still leads to large outliers (some of 8 orders of magnitude for 100 draws from the distribution), which are unphysical and continue to present difficulties when sampling with HMC.

To limit the number of unphysical outliers, we can instead limit the prior to consider only distributions that are less heavy-tailed than the Cauchy distribution — i.e. all those with $\nu > 1$. This is the approach taken by Gelman et al. (2013), rendering regions of parameter space inaccessible (in the case of the Gelman et al. (2013) prior, $\nu = 1$ is the cutoff, which corresponds to a Cauchy distribution.) On the other hand, the priors used in Juárez & Steel (2010); Ding (2014); Feeney et al. (2018) have disproportionate prior density at low values of $\nu$, which corresponds to models with outliers several orders of magnitude larger than the predicted effect size.

In this work, we use the prior

$$\nu \sim \text{Inv-}\Gamma(4, 15), \tag{A2}$$

where Inv-$\Gamma(\alpha, \beta)$ is an inverse gamma distribution with shape parameter, $\alpha$, and scale parameter, $\beta$.

This prior was chosen for two reasons:

• the prior is smooth in $\nu$ with no sharp boundaries (unlike Gelman et al. (2013))

• the prior has density at a larger range of outlier fractions $\omega$ than that in Juárez & Steel (2010) but insignificant density at unrealistically high outlier fractions, in contrast with the priors in Ding (2014); Feeney et al. (2018).

## APPENDIX B: ESTIMATING VARIANCE TRADE-OFFS

In Section 2.5, we examined the increased standard deviation in parameter estimates using a toy model, showing the results of repeated trials alongside the Cramér-Rao bound (Rao 1945; Cramér 1946) in Figure 4; here, we derive the formula for this bound.
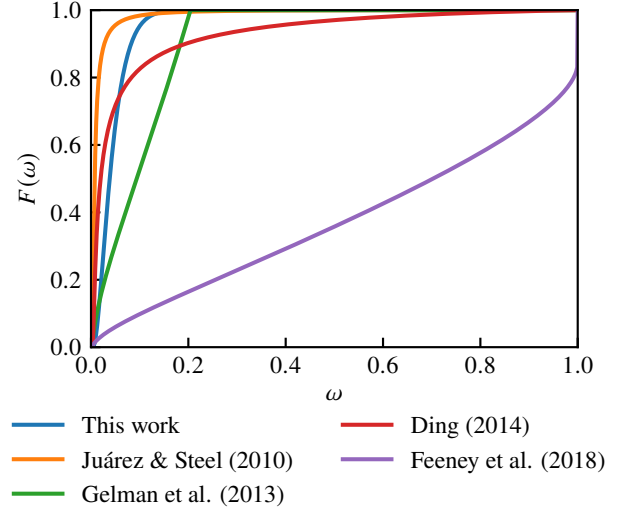
For the toy model, we looked at the posterior distribution of the mean $\mu$ for a series of $N$ points drawn from a normal distribution with zero mean and unit variance. The log likelihood, $\ell(\mu; \nu, \{x_i\})$, when analyzed with a $t$-distribution, is

$$\ell(\mu; \nu, \{x_i\}) = -\sum_i \frac{\nu + 1}{2} \log\left(1 + \frac{(x_i - \mu)^2}{\nu}\right) + C, \tag{B1}$$

where $C$ is a constant. Thus, the Fisher information is

$$\mathcal{F}(\mu; \nu, \{x_i\}) = (\nu + 1) \sum_i \frac{\nu - (x_i - \mu)^2}{\nu + (x_i - \mu)^2}. \tag{B2}$$

We can then evaluate the Cramér-Rao bound at the true mean $\mu = 0$, marginalising over the particular dataset $\{x_i\}$, to derive the variance in the estimate

$$\text{Var}(\hat{\mu}) \geq \frac{1}{N(\nu + 1)\left(1 - \sqrt{\frac{\nu}{2}\pi} \exp\left(\frac{\nu}{2}\right) \text{erfc}\left(\sqrt{\frac{\nu}{2}}\right)\right)}, \tag{B3}$$

where $\text{erfc}(x) = 1 - 2/\sqrt{\pi} \int_0^x \exp(-t^2)\, dt$ is the complementary error function. This bound is used to derive the posterior standard deviation ratio plotted in Figure 4.

## APPENDIX C: FIXED-VALUE CALIBRATION DATA MODELS

Here, we fully specify the dataset models used for fixed-value calibration tests in Section 4.

## C1  *t*-distributed data

In Section 4.1, the fixed-value test datasets have $N = 20$ datapoints drawn from the following distribution:

$$x_i \sim \mathcal{N}(\mu = 2, \sigma^2 = 4) \tag{C1}$$

$$y_i \sim t_3(\mu = 3 + 2x_i, \sigma^2 = 0.01/\sigma_{68}^2(\nu = 3)) \tag{C2}$$

$$\log_{10} \sigma_{x,i} \sim \mathcal{N}(\mu = -1, \sigma^2 = 0.01) \tag{C3}$$

$$\log_{10} \sigma_{y,i} \sim \mathcal{N}(\mu = -0.7, \sigma^2 = 0.01) \tag{C4}$$

$$\hat{x}_i \sim \mathcal{N}(\mu = x_i, \sigma^2 = \sigma_{x,i}^2) \tag{C5}$$

$$\hat{y}_i \sim \mathcal{N}(\mu = y_i, \sigma^2 = \sigma_{y,i}^2). \tag{C6}$$

## C2  Normally-distributed data with an outlier

In Section 4.2, we generated datasets of $N = 12$ points using the model:

$$x_i \sim \mathcal{N}(\mu = 5, \sigma^2 = 9) \tag{C7}$$

$$y_i \sim \begin{cases} \mathcal{N}(\mu = 3 + 2x_i - 10, \sigma^2 = 0.04) & \text{for the second-largest } x_i \\ \mathcal{N}(\mu = 3 + 2x_i, \sigma^2 = 0.04) & \text{otherwise} \end{cases} \tag{C8}$$

$$\log_{10} \sigma_{x,i} \sim \mathcal{N}(\mu = -0.5, \sigma^2 = 0.01) \tag{C9}$$

$$\log_{10} \sigma_{y,i} \sim \mathcal{N}(\mu = -0.3, \sigma^2 = 0.01) \tag{C10}$$

$$\hat{x}_i \sim \mathcal{N}(\mu = x_i, \sigma^2 = \sigma_{x,i}^2) \tag{C11}$$

$$\hat{y}_i \sim \mathcal{N}(\mu = y_i, \sigma^2 = \sigma_{y,i}^2). \tag{C12}$$

## C3  Two-dimensional normal mixture model

We introduce the parameter $O_i$ to indicate whether the $i$th datapoint is drawn from the core distribution (in which case, $O_i = 0$) or from the outlier distribution (for which $O_i = 1$).

In Section 4.3, we generated a dataset of $N = 200$ points using the model:

$$\boldsymbol{x}_i \sim \begin{cases} \mathcal{N}\left(\mu = \begin{pmatrix} -3 \\ 2 \end{pmatrix}, \Sigma^2 = \begin{pmatrix} 0.5 & -1 \\ -1 & 4 \end{pmatrix}^2\right) & 1 \leqslant i \leqslant 140 \\ \mathcal{N}\left(\mu = \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \Sigma^2 = \begin{pmatrix} 1 & 0.2 \\ 0.2 & 0.8 \end{pmatrix}^2\right) & 140 < i \leqslant 200 \end{cases} \tag{C13}$$

$$O_i \sim \text{Bernoulli}(0.1) \tag{C14}$$

$$y_i \sim \begin{cases} \mathcal{N}(\mu = 2 + (3, 1)^T \cdot x_i, \sigma^2 = 0.16) & O_i = 0 \\ \mathcal{N}(\mu = 2 + (3, 1)^T \cdot x_i, \sigma^2 = 16.0) & O_i = 1 \end{cases} \tag{C15}$$

$$\Sigma_{x,i} \sim \mathcal{W}_2(V = 0.1\mathbb{I}, n = 3) \tag{C16}$$

$$\log_{10} \sigma_{y,i} \sim \mathcal{N}(\mu = -1, \sigma^2 = 0.01) \tag{C17}$$

$$\hat{x}_i \sim \mathcal{N}(\mu = \boldsymbol{x}_i, \Sigma^2 = \Sigma_{x,i}^2) \tag{C18}$$

$$\hat{y}_i \sim \mathcal{N}(\mu = y_i, \sigma^2 = \sigma_{y,i}^2), \tag{C19}$$

where $\mathcal{W}_2$ denotes a Wishart distribution over 2x2 matrices and $\mathbb{I}$ denotes the 2x2 identity matrix.

## C4  Laplace-distributed data

In Section 4.4, we generate $N = 25$ datapoints under the following model:

$$x_i \quad \sim \quad \mathcal{U}(-5, 5) \tag{C20}$$

$$y_i \quad \sim \quad \text{Laplace}(\mu = -1 + 0.8x_i, b = 0.2) \tag{C21}$$

$$\log_{10} \sigma_{x,i} \quad \sim \quad \mathcal{N}(\mu = -1, \sigma^2 = 0.01) \tag{C22}$$

$$\log_{10} \sigma_{y,i} \quad \sim \quad \mathcal{N}(\mu = -1, \sigma^2 = 0.01) \tag{C23}$$

$$\hat{x}_i \quad \sim \quad \mathcal{N}(\mu = x_i, \sigma^2 = \sigma_{x,i}^2) \tag{C24}$$

$$\hat{y}_i \quad \sim \quad \mathcal{N}(\mu = y_i, \sigma^2 = \sigma_{y,i}^2). \tag{C25}$$

This paper has been typeset from a TeX/LaTeX file prepared by the author.