Gaussian Mixture Models for Blended Photometric Redshifts

Daniel M. Jones, 1 * Alan F. Heavens, 1

¹ Astrophysics Group & Imperial Centre for Inference and Cosmology, Imperial College London, London SW7 2A

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

Future cosmological galaxy surveys such as the Large Synoptic Survey Telescope (LSST) will photometrically observe very large numbers of galaxies. Without spectroscopy, the redshifts required for the analysis of these data will need to be inferred using photometric redshift techniques that are scalable to large sample sizes. The high number density of sources will also mean that around half are blended. We present a Bayesian photometric redshift method for blended sources that uses Gaussian mixture models to learn the joint flux-redshift distribution from a set of unblended training galaxies, and Bayesian model comparison to infer the number of galaxies comprising a blended source. The use of Gaussian mixture models renders both of these applications computationally efficient and therefore suitable for upcoming galaxy surveys.

Key words: cosmology: observations – galaxies: distances and redshifts – methods: statistical

of interest.

fluxes.

INTRODUCTION

Photometric galaxy surveys such as the VISTA Kilo-degree Infrared Galaxy (VIKING) survey (Edge et al. 2013) and the Dark Energy Survey (DES) (Dark Energy Survey Collaboration et al. 2016) have become important probes within current observational cosmology. These surveys use photometric observations of large samples of galaxies to probe the distribution of matter in the large scale structure of the Universe. This distribution is sensitive to several phenomena of interest to cosmology such as dark energy (e.g., DES Collaboration et al. 2018b), the rate of expansion described by the Hubble constant (e.g., Abbott et al. 2018), models beyond the standard flat Λ CDM model (e.g., DES Collaboration et al. 2018a) and the sum of the neutrino masses (e.g., Choudhury & Choubey 2018).

Making inferences about these phenomena requires the distribution of redshifts of galaxies in the sample. Spectroscopic observations that reach a sufficient signal-to-noise provide a way to obtain very precise redshifts. However, the size and depth of these galaxy samples render spectroscopy prohibitively time-consuming. As a result, photometric redshifts are a vital part of the analysis of cosmological galaxy

Photometric redshift methods can broadly be characterised into two types; template-based and machine learning methods. Template-based methods parametrise the relation

between flux and redshift through a set of spectral templates. Galaxy fluxes are forward modelled by redshifting

these spectra and integrating over the survey filters, allow-

ing the redshift to be inferred through standard maximum

likelihood (e.g., Bolzonella et al. 2000; Ilbert et al. 2006) or

Bayesian techniques (e.g., Benítez 2000). These template-

based methods are easily interpretable, and Bayesian infer-

ence allows rigorous statistical uncertainties to be propa-

gated through probability density functions (PDFs). How-

ever, their accuracy is dependent on the applicability of the

template sets, which are often small¹, to the galaxy sample

relation between flux and redshift from a training set of

galaxies with known redshifts. This relation is represented by

a flexible model such as random forests (e.g., Carliles et al.

2010; Carrasco Kind & Brunner 2013), boosted decision

trees (e.g., Gerdes et al. 2010), neural networks (e.g., Col-

lister & Lahav 2004; Sadeh et al. 2016), support vector ma-

chines (e.g., Wadadekar 2005) and Gaussian processes (e.g.,

Way & Srivastava 2006; Almosallam et al. 2016). Machine

learning approaches can also be extended to include extra

input features such as morphology (e.g., Soo et al. 2018) or to use entire images as input (e.g., D'Isanto & Polsterer

2018), rather than reducing this information to a vector of

Machine learning methods, on the other hand, learn the

¹ The photometric redshift software BPZ (Benítez 2000) is packaged with a set of 8 templates by default.

^{*} E-mail: d.jones15@imperial.ac.uk

The data-driven approach of machine learning methods avoids the potential pitfalls of small template sets, but instead relies on the training set being representative. If this is the case, the accuracy of these methods can be greater than that of than template-based methods (Hildebrandt et al. 2010) In practice however, training sets are often shallower than the photometric sample, reducing the accuracy of machine learning methods (Rivera et al. 2018).

The unrepresentativeness of training sets could be a problem for future surveys such as the Large Synoptic Survey Telescope (LSST; Ivezić et al. 2019) because their photometry will reach depths beyond which spectroscopy can be reasonably performed. However, by increasing the number density of galaxies on the sky, these very deep observations enable cosmological constraints with higher precision than current galaxy surveys.

Another major challenge presented by the depth of future surveys like LSST is blending (Chang et al. 2013), the chance overlapping of galaxies along the line of sight. As a result, the intrinsic fluxes of blended galaxies cannot be observed directly, only their noisy combination. The focus of this paper is to address the challenge of inferring redshifts from this blended photometry in a way that can scale to the large datasets of future surveys.

One approach to counter the problem of blending is to separate blended sources into distinct images of each constituent galaxy, known as deblending. Deblending methods that rely solely on the morphological information contained within a single band (e.g., Lupton 2005) will separate blended galaxies with large angular separations more easily than those that are more closely aligned. Dawson & Schneider (2014) found that for a survey like LSST where 44-55%of sources are blended, these methods would misidentify 15-20% of all sources as unblended. As a result, more recent deblending methods (e.g., Joseph et al. 2016; Melchior et al. 2018) also utilise colour information by using galaxy images in several bands.

When splitting the analysis into separate images in this way, it is important to take care with how uncertainties from the deblending process are propagated. While the total flux of a source may be well constrained by observations, the separate flux of each galaxy is not, as it is not observed independently of other galaxies it is blended with. As a result, the errors on the fluxes of each galaxy will be correlated. Ideally, this correlation should be propagated to later analyses, though these uncertainties can be difficult to estimate and propagate for these deblending methods (see, e.g., Melchior et al. 2018).

An alternative to deblending is to infer quantities of interest, such as photometric redshifts, from blended data directly. This joint approach automatically accounts for correlations between each galaxy in a blended source and correctly propagates these uncertainties to the final results. This is the approach taken in Jones & Heavens (2019), which generalises Bayesian template-based photometric redshift methods to the case of blended observations.

This paper takes the same joint-inference approach, but uses a Gaussian mixture model to learn the flux-redshift relation from a training set of galaxies with known redshifts. We then use this model as a prior to derive the posteriors and marginal-likelihoods for sources consisting of one or two galaxies. Since these can be computed analytically,

Table 1. A summary of the notation used throughout this paper.

Symbol	Description
N	Number of constituent galaxies in a source
z_n	Model redshift of constituent galaxy n
F_n	Model flux vector of constituent galaxy n
$\overline{\hat{F}}^n$	Vector of observed fluxes
$rac{ar{F}}{\hat{F}}^n \ \Sigma^{\hat{F}}$	Covariance matrix of observed fluxes
M	Number of components in the mixture model
w^k	Weight of mixture component k
$rac{\mu^k}{\Sigma^k}$	Mean vector of mixture component k
$\sum_{k}^{\infty} k$	Covariance matrix of mixture component k
\mathcal{E}^1	Evidence for single-constituent model
\mathcal{E}^2	Evidence for two-constituent model
$\mathcal{N}(x \mid \mu, \Sigma)$	Multivariate Gaussian PDF with mean vector
	μ and covariance matrix Σ
$\mathcal{N}(\underline{x} \mid \underline{\mu}, \Sigma)$ $\tilde{\mathcal{N}}(\underline{x} \mid \underline{\eta}, \Lambda)$	Multivariate Gaussian PDF in natural
	parametrisation with parameters $\Lambda \equiv \Sigma^{-1}$ and
	$\underline{\eta} \equiv \Sigma^{-1} \underline{\mu}$

this is significantly less computationally demanding than the nested sampling-based method described in Jones & Heavens (2019), an important property for use in future galaxy surveys.

Gaussian mixture models have previously been used for obtaining photometric redshifts of quasars Bovy et al. (2012). The method presented in this paper is an extension of this approach to deriving posterior distributions and Bayesian evidences for the redshifts of blended sources. This also builds on our previous work (Jones & Heavens 2019) by being completely data-driven, learning the mapping between flux and redshift from a training set, rather than imposing it a priori through a set of templates. The significant computational advantages afforded by this approach now allows blended photometric redshifts to be applied to very large future datasets.

Throughout this paper, we use the term *constituent* to describe the individual galaxies comprising a blended source. Following convention, we refer to each multivariate Gaussian distribution in the mixture model as a component. We denote scalars using an italic font x, vectors using an underlined italic font x and matrices using a bold italic font x. We summarise our notation in Table 1.

In section 2, we briefly describe the results of Jones & Heavens (2019) and introduce our blended photo-z formalism. In section 3, we introduce our formalism for blended photometric redshifts with Gaussian mixture models. We use this to derive expressions for the posteriors and evidences in section 4. We present results of tests of our method on simulated data in section 5. Finally, in section 6, we present these tests on real blended data from the Galaxy And Mass Assembly (GAMA) survey (Baldry et al. 2017).

BLENDED PHOTOMETRIC REDSHIFTS

A Bayesian method for inferring the photometric redshifts of blended sources was introduced in Jones & Heavens (2019). This is a template-based method, generalising the commonly used Bayesian Photo-z (BPZ) method of Benítez (2000). A summary of the main result, the joint posterior distribution of the redshift and magnitude of each galaxy within a blended source, is given below.

For a given template t at redshift z, the model flux $T_{t,b}(z)$ in each band b can be calculated by integrating the redshifted galaxy spectrum over the respective filter response. These fluxes are then scaled such that the flux in an arbitrarily chosen reference band b_0 is equal to $10^{-0.4m_0}$, where the magnitude m_0 is a parameter to be inferred in addition to the redshift. The predicted flux for a blended source is then given as a linear combination of these galaxy fluxes, i.e.,

$$F_{\{t\},b}^{(N)}\left(\{z\},\{m_0\}\right) = \sum_{\alpha=1}^{N} \frac{10^{-0.4m_{0,\alpha}}}{T_{t_\alpha,b_0}(z_\alpha)} T_{t_\alpha,b}\left(z_\alpha\right), \tag{1}$$

where z_{α} , $m_{0,\alpha}$ and t_{α} are the redshift, reference-band magnitude and template for constituent α respectively, and N is the number of constituents in the source. The desired posterior can then be found by marginalising over the template for each galaxy and applying Bayes rule to give

$$P(\{z\}, \{m_0\} \mid \hat{\underline{F}}, N) \propto \sum_{i=1}^{T^N} P(\hat{\underline{F}}, \mid \{z\}, \{t\}_i, \{m_0\}, N) P(\{z\}, \{t\}_i, \{m_0\} \mid N).$$
(2)

where $\hat{\underline{F}}$ is the vector of observed fluxes and T is the number of templates in the template set. The predicted flux in equation 1 is defined for a particular choice of template for each galaxy within the blended source. The template marginalisation therefore runs over the T^N combinations of this choice. Note that this posterior is conditioned on a particular choice of N, the number of galaxies within the blended source; setting this is described in section 2.1.

The joint prior can then be developed by factorising into priors defined for each constituent. In doing this, three blending-specific complications arise. Firstly, the redshifts of each constituent are not independent since galaxies are clustered. As a result, the posterior for N blended sources should include an additional term involving correlation functions up to N-point to account for this.

Secondly, the effect of source selection should also be accounted for. One effect of this is that the selection criteria imposes a faint-end cut on the magnitude prior. Without this cut, its simple analytic form would be improper, rendering the model selection described in section 2.1 impossible.

Lastly, a sorting condition is required to break the exchangeability of the blended constituents. Allowing this exchangeability results in marginal redshift distributions with contributions from every constituent, i.e., they would always have multiple peaks. By enforcing an ordering, the posterior better represents the underlying physical source. Jones & Heavens (2019) found that redshifts were recovered more successfully when applying this sorting condition to constituent redshifts, though sorting the magnitudes is also sufficient to break the exchangeability.

2.1 Model selection for identifying blends

In addition to inferring the redshift of each constituent in a blended source, the method of Jones & Heavens (2019) can also identify whether a source is blended. Since the posterior defined in equation 2 is conditioned on the number of

constituents N, we can consider this choice to be the model and use Bayesian model comparison techniques to infer the number of constituents within the source.

To compare two models with a source of n and m constituents, we write the relative probability and apply Bayes rule to give

$$\mathcal{P}_{n,m} = \frac{P(N=n \mid \underline{\hat{F}}, \hat{F}_0)}{P(N=m \mid \underline{\hat{F}}, \hat{F}_0)} = \frac{P(N=n)}{P(N=m)} \frac{P(\underline{\hat{F}}, \hat{F}_0 \mid N=n)}{P(\underline{\hat{F}}, \hat{F}_0 \mid N=m)}.$$
(3)

The first term is the ratio of model priors, allowing the $a\ priori$ probability of a source being blended to be set. This value could be informed by the expected number of blended sources given the survey depth, or could leverage additional independent information such as whether a source is located within a cluster or the field. Throughout, we assume this ratio is unity, so that a source is equally likely to be blended as not. This assumption is trivial to modify, however.

The second term in equation 3 is a ratio of marginal likelihoods known as the Bayes factor. The calculation of these marginal likelihoods, also known as evidences and labelled \mathcal{E} , involves an integral over the full support of the prior, i.e.,

$$\mathcal{E} \equiv P(\underline{d} \mid \mathcal{M}) = \int P(\underline{d} \mid \{\theta\}, \mathcal{M}) P(\{\theta\} \mid \mathcal{M}) d\{\theta\}, \tag{4}$$

where \underline{d} is the data vector, \mathcal{M} is the model and the integral is over the set of model parameters $\{\theta\}$.

This integral is often difficult to evaluate, particularly if the dimensionality of the parameter space is large. While the prior volume may be large, the likelihood can peak sharply. Nevertheless, the comparatively low-density tails of the posterior can contain significant volume and can therefore not be ignored. Numerically evaluating an integral with nonnegligible contributions at both of these scales is computationally challenging.

In order to sample the posterior and evaluate the evidence, Jones & Heavens (2019) uses MultiNest (Feroz et al. 2009), an efficient implementation of the nested sampling method (Skilling 2006). However, even sampling with an efficient method such as MultiNest can be computationally demanding; sampling both the two- and four-dimensional posteriors for one- and two-constituent sources respectively² takes approximately two minutes per source on a workstation with a 3 GHz Intel Xeon processor. While this is viable for small samples, it is not scalable to the large samples of $\sim 10^9$ galaxies in a future survey like LSST. Instead, this paper develops a method that does not rely on these computationally demanding integrals. As a result, the one- and two-constituent inference and model selection can be done for approximately ten sources per second, a speed-up of three orders of magnitude on the workstation described above. Photometric redshift inference is also trivially parallelisable for high-performance computing environments, since each source can be considered independently. We present our method below.

 $^{^2\,}$ Both a redshift and a magnitude is sampled for each constituent in the source.

3 GAUSSIAN MIXTURE MODEL PHOTO-Z

Photometric redshifts inferred using machine learning methods are often very accurate when good training data is available. These methods perform regression, and use this training data to learn the mapping from fluxes to redshifts. Many machine learning algorithms are not inherently probabilistic; a particular input will map to a particular output. However, accurate uncertainties on cosmological parameters rely on propagating uncertainties from all stages of the analysis. Machine learning photometric redshift methods have therefore developed several ways to estimate these uncertainties.

One example that accounts for errors in the observed fluxes is to apply the chain rule to successive layers of a neural network (Collister & Lahav 2004), providing the variance of the output redshift. Some machine learning methods such as a Gaussian process (e.g., Way & Srivastava 2006), are already explicitly probabilistic, naturally producing variance estimates alongside their prediction. Other methods can represent their uncertainties more generally by inferring PDFs as their output. This can be done by training many machine learning algorithms to each independently estimate the redshift and taking the distribution of the ensemble to be the redshift PDF (Sadeh et al. 2016). A single neural network can also accomplish this by being trained to output the parameters of a parametrised PDF rather than the redshift directly (D'Isanto & Polsterer 2018). PDFs represent the complete probabilistic knowledge over a system under investigation, and are thus a general mechanism for quantifying and propagating uncertainties within a statistical analysis (e.g. Gelman et al. 2013).

In addition to enabling the rigorous propagation of uncertainties, using full photometric redshift PDFs has been shown to improve the accuracy of cosmological analyses (e.g., Mandelbaum et al. 2008; Myers et al. 2009). PDFs also have an advantage over simply representing uncertainty with the variance in their ability to represent multimodality; that is, several distinct, well separated redshifts being plausible for a given vector of fluxes. This is a common occurrence in photometric redshifts (Benítez et al. 2009). Colour-redshift degeneracies mean that high- and low- redshift galaxies can have similar colours, often due to spectral features such as the Lyman and Balmer breaks being misidentified as one another (Graham et al. 2018).

Here, we treat the training data not as variables to regress between, but instead as noisy samples from the joint redshift-flux distribution, turning the problem into one of density estimation. The joint density is the most general probabilistic description of the training data, allowing several quantities of interest to be derived. Given an observed vector of fluxes \hat{F} , the redshift can be inferred using the conditional distribution $P(z \mid \hat{F})$ which can be derived from the joint distribution. This PDF can be multimodal, capturing the degeneracy described above. These distributions can be composed together to produce the conditional distribution of the redshifts of a blended source $P(z_1, z_2 \mid \hat{\underline{F}})$ in a similar fashion. The joint distribution also permits calculation of marginal likelihoods, allowing Bayesian model selection techniques to be used to infer the number of constituents in a source. Finally, the interpretation of the joint distribution is clear, in contrast to other machine learning methods that can be 'black-boxes', requiring additional ad-hoc techniques to improve their interpretability (e.g., Shrikumar et al. 2017; Shwartz-Ziv & Tishby 2017)

We model the joint distribution of the latent, noise-free parameters as a Gaussian mixture model (GMM), a weighted linear combination of multivariate Gaussians, i.e.,

$$P(z, \underline{F}) = \sum_{k} w^{k} \mathcal{N}(z, \underline{F} \mid \underline{\mu}^{k}, \Sigma^{k}).$$
 (5)

By imposing that $\sum_k w^k = 1$, this density is correctly normalised, i.e.,

$$\sum_{k} w^{k} \iint \mathcal{N}(z, \underline{F} \mid \underline{\mu}^{k}, \Sigma^{k}) \, dz \, d\underline{F} = \sum_{k} w^{k} = 1.$$
 (6)

This choice has several useful features. Firstly, GMMs are easy to train using standard, well-tested methods. This is discussed further in section 3.1. Secondly, inference with GMMs is computationally inexpensive as they can be efficiently sampled as detailed in section 3.4. Lastly, GMMs are mathematically convenient. Both the conditional and marginal distributions of multivariate Gaussians are also Gaussians. The same is also true of both the product and convolution of several multivariate Gaussians. These properties will be used frequently throughout this paper to render many calculations analytic. Despite this, GMMs can represent a wide variety of PDFs, including those that are skewed or multimodal. This is demonstrated in Fig. 1.

Using GMMs to infer photometric redshifts in this way was first done in Bovy et al. (2012), who applied the method to obtain photometric redshifts of quasars and used model selection techniques to separate stars and quasars. The method we present in this paper differs from this in several ways. Firstly, we extend the method to the case of jointly inferring multiple redshifts directly from blended data.

Secondly, Bovy et al. (2012) fit a series of many GMMs to the fluxes and redshifts of quasars in several magnitude bins. As a result, our model has significantly fewer parameters to fit. Nevertheless, our use of cross-validation to set the number of mixture components as described in section 3.3 provides the model sufficient flexibility to fit the flux-redshift density with the full fidelity provided by the training set.

The binning of Bovy et al. (2012) is not possible due to the extension to blended sources. Observations in this case are of the flux of the blended source, while the magnitude bin in that model is chosen based on the magnitude of an individual galaxy. This quantity that is not observed in the blended case, and so cannot be used to choose a magnitude bin. The same is true of colours, i.e., ratios of fluxes relative to the flux in a particular reference band, which are often used in machine learning-based photometric redshift methods. Since the reference-band flux of each galaxy in a blended source is not observed, the colours for each galaxy cannot be calculated and so cannot be used to infer the redshifts.

Finally, our derivation does not use the convolution property of multivariate Gaussians described above, since integrals over fluxes are then implicitly evaluated from $-\infty$ to ∞ as multivariate Gaussians have infinite support. These integrals therefore contain contributions from non-physical negative fluxes. This is a safe approximation when considering unblended sources, since their flux is strongly constrained by observations. However, the same is not true of blended sources, where the individual flux of each constituent is not observed. Instead, we evaluate these results

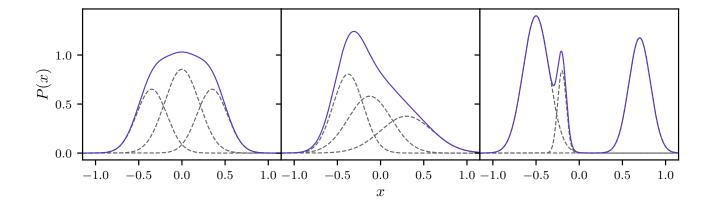


Figure 1. Plot showing a variety of PDFs that can be represented by Gaussian mixture models, given a sufficient number of components. The dashed grey curves show each weighted Gaussian component, and the solid blue curves show the mixture formed by the linear combination of these components.

using an efficient Monte Carlo integration method. We therefore treat unblended sources in the same way for consistency.

All fluxes throughout are renormalised for numerical stability. This is done by dividing each flux by the standard deviation in the training set, e.g., for band b,

$$F_b \to \frac{F_b}{\sigma(\{\hat{\underline{F}}_b\}_{\rm tr})}$$
 (7)

Normalising the data in this way is a common preprocessing step in machine learning methods. Without this renormalisation, the observed fluxes are small enough that the EM fitting procedure is dominated by numerical errors as the covariance matrices of the components become poorly conditioned. The corresponding change in the covariance matrix of each data point is given by

$$\Sigma_{ij} \to \frac{\Sigma_{ij}}{\sigma(\{\underline{\hat{F}}_i\}_{\mathrm{tr}})\sigma(\{\underline{\hat{F}}_j\}_{\mathrm{tr}})}$$
 (8)

We also note that magnitudes are commonly used for this purpose in machine learning-based photometric redshift methods, since the logarithmic transformation of the flux also effectively normalises them. However, an advantage of the GMM method presented here is that expressions for posteriors and evidences can be calculated analytically. This relies on the model for the flux of the blended sources being a linear combination of the fluxes of the individual constituents, since this leaves the likelihood of the sum a Gaussian. This would no longer be the case when using magnitudes, as the model for the magnitude of a blended source would be a non-linear function of the individual constituent magnitudes.

3.1 Training Gaussian mixture models

Our prior density $P(z, \underline{F})$ is defined in terms of the true, latent parameters. Therefore, this density must be fitted with a method that incorporates both the noisy data and the covariance. To do this, we use extreme deconvolution (Bovy et al. 2011), an extension of the expectation-maximisation (EM) algorithm (Dempster et al. 1977) commonly used the

find the maximum-likelihood parameters of GMMs. This is the same fitting method as the quasar photometric redshift method of Bovy et al. (2012).

Extreme deconvolution generalises the EM algorithm to the case where the data is subject to normally-distributed errors. The EM algorithm is a general method for fitting models with some form of hidden data in addition to the observed data. Given an initial guess at the parameters, the algorithm iteratively modifies these parameters to increase the likelihood, converging to a local maximum.

For a single multivariate Gaussian, the maximum-likelihood parameters can be found exactly through the derivative of the likelihood. However, the same is not true of mixtures of Gaussians, as these parameters are not available in closed form. The hidden information that would make this tractable is the identity of the component from which each sample was drawn. If this were known, fitting the GMM would reduce to the previous analytic case. Though this information is hidden, this points to an iterative solution; first, the parameter guess can be used to update the hidden information, then this information can be used to update the parameters.

In essence, expectation-maximisation is a probabilistic version of this procedure that takes into account the uncertainty in the hidden information. By averaging the likelihood over the probability of each sample being drawn from each component, the maximum likelihood parameters can be found in closed form. Since the component probability depends on the parameters being fitted, this process is iterative.

The extreme deconvolution method of Bovy et al. (2011) extends the EM algorithm to fit data with Gaussian errors. This is done by replacing the likelihood with a marginalised version given by

$$P(\underline{\hat{x}} \mid \{\theta\}) = \int P(\underline{\hat{x}}, \underline{x} \mid \{\theta\}) \, d\underline{x} = \int P(\underline{\hat{x}} \mid \underline{x}) P(\underline{x} \mid \{\theta\}) \, d\underline{x},$$

$$(9)$$

where \hat{x} is the vector of observed values, \underline{x} is the latent vector of true values and $\{\theta\}$ are the mixture parameters being fitted, i.e., weights, means and covariances. The data likeli-

hood $P(\underline{\hat{x}} \mid \underline{x})$ is assumed to be a multivariate Gaussian, and $P(\underline{x} \mid \{\theta\})$ is the GMM. Due to the convolution property of multivariate Gaussians, this marginalised likelihood is also a Gaussian mixture, and thus amenable to being fitted using an expectation-maximisation approach. Using this extreme deconvolution method, we fit the joint flux-redshift distribution $P(z,\underline{F})$ while accounting for uncertainties in the training set.

This fitting procedure assumes that the number of mixture components is fixed. The method we use to decide on this number is discussed in section 3.3.

As discussed above, multivariate Gaussians have infinite support, and so non-physical negative fluxes and negative redshifts are a priori allowed. No non-physical fluxes will be present in the training set, and negative redshifts, while not non-physical, are sufficiently rare that they can be presumed to not be present either. As a result, there is no incentive for the training algorithm to assign significant prior volume here. However, without an additional prior on the mixture parameters, prior volume in negative regions is not penalised either.

It is possible to generalise EM-based methods such as extreme deconvolution to maximise the posterior rather than the likelihood by adding a log-prior. However, while this will ameliorate the problem of negative values, it cannot eliminate it completely; the GMM having infinite support means that every point in parameter space will always have non-zero density.

An alternative approach is to impose an additional prior that is zero is any negative regions of parameter space, i.e.,

$$P(z, \underline{F}) = \Theta(z, \underline{F}) \sum_{k} w^{k} \mathcal{N}(z, \underline{F} \mid \underline{\mu}^{k}, \Sigma^{k})$$
 (10)

where

6

$$\Theta(z, \underline{F}) = \begin{cases} 0 & \text{for } z, \underline{F} < 0 \\ 1 & \text{otherwise.} \end{cases}$$
 (11)

This will exactly fix the problem of negative values. However, it will also force otherwise analytic integrations to have to be done numerically. These cases are discussed in the relevant sections below.

Imposing this boundary prior will also change the normalisation of the prior from unity, i.e.,

$$\iint \Theta(z,\underline{F}) \sum_{k} w^{k} \, \mathcal{N}(z,\underline{F} \mid \underline{\mu}^{k}, \Sigma^{k}) \, \mathrm{d}z \, \mathrm{d}\underline{F} \neq 1 \, . \tag{12}$$

The model selection described below requires that the prior be normalised. This normalisation differs between the singleand two-constituent cases, with the latter also being affected by the sorting condition. These normalisations are therefore discussed in their respective sections below.

It should be noted that, since this is an empirical method that does not rely on any underlying physical model in the way that a template-based method does, the redshift can be transformed almost arbitrarily. The only restrictions in this transformation are that it is both invertible and well-defined for all positive real numbers. The only modifications to the method required to accommodate this are to the limits of redshift integrals. For a transformation $\mathcal{T}(z)$, the lower and upper limits should be replaced with $\mathcal{T}(0)$ and $\mathcal{T}(\infty)$ respectively.

The transformation $\mathcal{T}(z) = \log(z)$ would seem to be a

sensible choice, as the lower and upper integration limits would become $-\infty$ and ∞ respectively, rendering all the redshift integrations throughout analytic. This is the approach taken by Bovy et al. (2012). However, in our tests, we found that this transformation reduces the accuracy of the blended redshift inference. The difference in accuracy of the single redshift inference was negligible. As a result, we do not transform redshifts throughout this paper.

A plot of this prior distribution, fitted to the simulated LSST-like training data described in section 5 and plotted using corner.py (Foreman-Mackey 2016), is shown in Fig. 2. The ability to plot this distribution is an advantage to this GMM method. As described above, machine learning methods can act as black boxes, where what has been learned is a complicated function approximator that can be difficult to interpret. In contrast, the central object being learned here is the joint flux-redshift distribution, a meaningful statistical object that can be plotted, sampled from and manipulated mathematically.

3.2 Utilising blended training data

The derivations detailed in sections 4.1 and 4.2 are presented for a scalar redshift z. However, it should be noted that these single-constituent results also hold for a vector \underline{z} . As a result, this method can be generalised so that the model is fitted to blended training data, i.e., a vector of blended fluxes with the associated vector of redshifts for each constituent.

Utilising blended training data would allow the method to infer both redshifts and the number of constituents accurately in cases where the blended constituents were systematically different from non-blended constituents. The cost of this, however, is an increase in the required size of the training set. Machine learning-based methods require a training set that is representative of the test set in order to be accurate. A blended training set would therefore have to contain sufficient examples of all possible pairs of constituents, rather than the constituents alone as required for the results in sections 4.3 and 4.4.

3.3 Cross-validating the number of mixture components

The procedure described in section 3.1 will fit the weights, means and covariances of the GMM for a fixed number of components. However, it is difficult a priori to choose this number; including more components within the mixture allows it more flexibility, but too many will cause the model to overfit. Given enough mixture components, the variance of each component will approach zero, with each being responsible for only a single sample. While this will significantly increase the likelihood of the training set, it will also cause the model to generalise extremely poorly.

Overfitting is a general concern when fitting machine learning models. As a result, various techniques for preventing overfitting have been suggested. These include restricting the dimensionality of the parameter space as we do here by fixing the number of components, disfavouring overfitted parameters through regularisation (e.g., Hoerl & Kennard 1970) or Bayesian priors (e.g., MacKay 1992), and stopping training before overfitting occurs (e.g., Prechelt 1998).

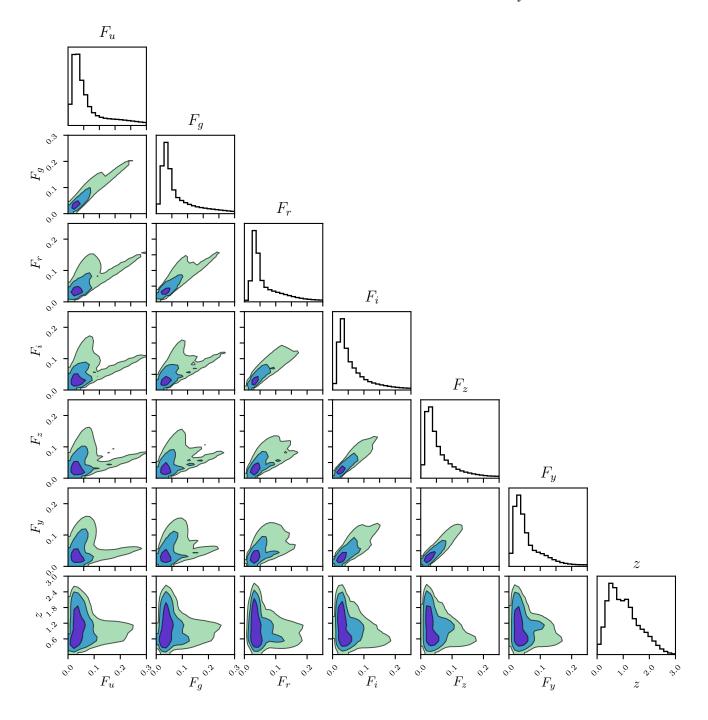


Figure 2. Corner plot of an example flux-redshift distribution fitted by our model. This density shown here is visualised using 10^6 samples drawn from a model that was fitted to the LSST-like simulations presented in section 5.

The ability for a machine learning method to generalise and whether it has been overfitted can be tested by using a a validation set, an additional set of data where the input and output are known but is not used during the training. By measuring the difference between the prediction and the known ground truth, the model can be evaluated.

It is useful to point out that a corollary to the notion of overfitting is that the fitting procedure need not converge to a global maximum, as that set of parameters will overfit the data. Instead, local maxima can be nearly as accurate on the test set, while generalising much better (Choromanska et al. 2014). Therefore, it is reasonable to use parameters corresponding to local maxima that are found to perform well during validation. This can avoid expending significant optimisation effort attempting to fit the global maximum.

To choose the number of components, we use k-fold cross validation, a method that repeatedly splits the data into training and validation sets. The training set is first split into k subsets. The model is then trained on k-1 subsets of this data, assuming a fixed number of mixture compo-

nents M. The remaining subset is then used for validation. By evaluating the model using the fluxes of this subset, the redshift predictions can be compared to the known truth and scored based on their accuracy. This training and validation is repeated k times for each number of components considered so that each subset is used for evaluation once. The average score can then be used to evaluate each number of components.

To evaluate the accuracy of the redshift predictions, we use the RMS scatter. Given a predicted redshift $z_{\mathrm{p},g}$ and a spectroscopic redshift $\hat{z}_{\mathrm{s},g}$ for galaxy g, the normalised error is defined as

$$\tilde{\delta}_g = \frac{\hat{z}_{s,g} - z_{p,g}}{1 + \hat{z}_{s,g}} \,. \tag{13}$$

After calculating this error for $n_{\rm g}$ galaxies, the RMS scatter for the sample is then given by

$$\sigma_{\rm RMS} = \sqrt{\frac{1}{n_{\rm g}} \sum_{g} \tilde{\delta}_{g}^{2}} \,. \tag{14}$$

This metric is evaluated using k-fold validation for each number of mixture components M being considered. We then choose M to be the number of components that minimises the RMS scatter averaged over each of the k folds.

3.4 Sampling from Gaussian mixture models

One of the significant advantages of using GMMs is that they can be efficiently sampled from without using methods such as MCMC. Since they are simply linear combinations of component distributions, a simple sampling scheme is to randomly select one of the components with a probability given by the weights, and then to draw a sample from the respective multivariate Gaussian.

This sampling scheme allows GMMs to be sampled efficiently and without rejection. However, the addition of the boundary prior described in section 4.1 means that samples with negative fluxes and redshifts are rejected during inference. Nevertheless, the efficiency of this sampling scheme means that this does not pose a problem, since many samples can still be drawn from the relevant posterior with little computational effort.

3.5 Compressed storage of PDFs

As described above, it is important that the results of photometric redshifts are represented as a PDF. However, given the large sample sizes of future galaxy surveys like LSST, storing these PDFs can present a problem. While a point estimate of the redshift and an associated error can be stored simply as two real numbers, a PDF will generally require many more. A naive representation of this distribution is a histogram where the redshift bins are fixed for all sources. While this is simple, it is not space efficient.

This problem was first investigated by Carrasco Kind & Brunner (2014), which proposed a sparse basis representation using Gaussian and Voigt distributions. Using this method, the PDF can be stored in a single signed integer per basis function, with O(10) basis functions required to accurately reconstruct the original PDFs. Malz et al. (2018) test PDF compression methods by measuring the Kullback-Leibler divergence between the original and compressed

PDFs. They suggest storing the redshifts corresponding to equally-spaced quantiles as an alternative to histograms.

The posteriors presented here are GMMs, potentially multiplied by an additional physical constraint. This representation permits a simple compression technique of discarding low-weight components. By construction, the number of components in the mixture describing the prior is the same as the mixture describing the redshift posterior. However, the latter is generally significantly more compact, describing the density over the parameter space for a single source only, rather than the entire population. It is therefore reasonable to expect that this posterior distribution could be represented by fewer components than the prior.

If additional computation can be afforded for a further reduction in storage space, mixture components can also be merged into a smaller number of approximating components. This procedure is known as mixture reduction (see, e.g., West 1993; Williams & Maybeck 2006; Runnalls 2007; Schieferdecker & Huber 2009).

4 DERIVING POSTERIORS AND EVIDENCES

4.1 Single-constituent posterior

We now derive the posterior distribution assuming that the source consists of a single, unblended constituent galaxy. The redshift under this model can then be inferred by sampling from this posterior, as described in section 3.4. We start by marginalising over the true, latent flux vector \underline{F} , giving

$$P(z \mid \underline{\hat{F}}) = \int P(z, \underline{F} \mid \underline{\hat{F}}) \, d\underline{F}.$$
 (15)

Applying Bayes rule, this becomes

$$P(z \mid \underline{\hat{F}}) \propto \int P(\underline{\hat{F}} \mid \underline{F}) P(z, \underline{F}) d\underline{F},$$
 (16)

where the unnecessary redshift conditioning has been dropped from the likelihood. We assume the likelihood to be a multivariate Gaussian centred on the observed fluxes, i.e.,

$$P(\underline{\hat{F}} \mid \underline{F}) = \mathcal{N}(\underline{F} \mid \underline{\hat{F}}, \Sigma^{\widehat{F}}), \tag{17}$$

where $\Sigma^{\hat{F}}$ is the covariance matrix of the observation. Galaxy surveys typically assume the errors on observed fluxes in each band to be independent, i.e., given as a flux and an error. In this case, the covariance matrix would simply be diagonal. No assumption is made about this covariance throughout however, allowing fluxes to be correlated in general.

The prior in equation 16 is given by the GMM described above. This prior is the only term involving the redshift; it fully represents the relation between flux and redshift learned from the training set.

Inserting both the prior and the likelihood into equation 16, the posterior becomes

$$P(z \mid \underline{\hat{F}}) \propto \sum_{k} w^{k} \int \mathcal{N}(\underline{F} \mid \underline{\hat{F}}, \Sigma^{\hat{F}}) \mathcal{N}(z, \underline{F} \mid \underline{\mu}^{k}, \Sigma^{k}) d\underline{F}.$$
 (18)

This posterior now contains the product of two Gaussian PDFs, albeit with different dimensionalities. We proceed by combining these two densities into a single multivariate Gaussian. This is analogous to the derivation of Bovy

et al. (2012). However, as described above, we do not make use of the convolution property of multivariate Gaussians, instead forming the product explicitly. To do this, we write our posterior in terms of a parameter vector $\underline{\theta}$ partitioned into redshift and fluxes, i.e.,

$$\underline{\theta} = \begin{pmatrix} z \\ F \end{pmatrix} . \tag{19}$$

Throughout, we label the redshift and flux blocks of parameters partitioned in the same way with z and f respectively.

The likelihood involves only the flux partition of the parameter vector. However, our prior has support over both redshift and flux, i.e., all of $\underline{\theta}$. The component parameters are thus partitioned in the same way so that the mean and covariance are given by

$$\underline{\mu}^k = \begin{pmatrix} \mu_z^k \\ \underline{\mu}_f^k \end{pmatrix} \tag{20}$$

and

$$\Sigma^{k} = \begin{pmatrix} \Sigma_{\text{ZZ}}^{k} & \frac{\Sigma_{\text{Zf}}^{k}}{\Sigma_{\text{ff}}^{k}} \\ \frac{\Sigma_{\text{fg}}^{k}}{\Sigma_{\text{ff}}^{k}} \end{pmatrix}$$
 (21)

respectively. The product of these two densities is most easily written in terms of the natural parametrisation³ of the multivariate Gaussian. This has a density given by

$$\tilde{\mathcal{N}}(\underline{x}|\underline{\eta}, \mathbf{\Lambda}) = \exp\left[\alpha + \underline{\eta}^T \underline{x} - \frac{1}{2}\underline{x}^T \mathbf{\Lambda}\underline{x}\right], \tag{22}$$

where we have added a tilde to notate the alternative parametrisation. The normalisation factor is given by

$$\alpha = -\frac{1}{2} \left[d \log(2\pi) - \log |\mathbf{\Lambda}| + \underline{\eta}^T \mathbf{\Lambda}^{-1} \underline{\eta} \right], \tag{23}$$

and the covariance matrix and mean vector are replaced with the natural parameters $\mathbf{\Lambda} \equiv \mathbf{\Sigma}^{-1}$ and $\underline{\eta} \equiv \mathbf{\Sigma}^{-1}\underline{\mu}$. respectively. The inverse covariance matrix $\mathbf{\Lambda}$ is known as the precision matrix. The product of the two densities in equation 18 can then be combined into a single multivariate Gaussian written in this natural parametrisation, given by

$$\mathcal{N}(\underline{F} \mid \underline{\hat{F}}, \Sigma^{\hat{F}}) \mathcal{N}(z, \underline{F} \mid \underline{\mu}^{k}, \Sigma^{k}) = c\tilde{\mathcal{N}}(z, \underline{F} \mid \underline{\eta}^{k\hat{F}}, \Lambda^{k\hat{F}}), \tag{24}$$

where the new parameters are

$$\mathbf{\Lambda}^{k\hat{F}} = (\mathbf{\Sigma}^k)^{-1} + \begin{pmatrix} 0 & 0 \\ 0 & (\mathbf{\Sigma}^{\hat{F}})^{-1} \end{pmatrix}$$
 (25)

and

$$\underline{\eta}^{k\hat{F}} = (\Sigma^k)^{-1}\underline{\mu}^k + \begin{pmatrix} 0 \\ (\Sigma^{\hat{F}})^{-1}\underline{\hat{F}} \end{pmatrix}. \tag{26}$$

Conveniently, the constant of proportionality c can also be written in terms of a multivariate Gaussian in standard parametrisation. This is given by

$$c_1^k = \mathcal{N}(\underline{\mu}_f^k \mid \underline{\hat{F}}, \Sigma_{ff}^k + \Sigma^{\hat{F}}). \tag{27}$$

These results are close to a standard property (e.g., Petersen & Pedersen 2014) where the product of two multivariate Gaussian densities is also a multivariate Gaussian. However,

the differing dimensionalities of the two densities in equation 24 slightly alter the expressions for the new parameters.

Inserting these results into equation 18 and moving constant terms outside of the integral, the expression for the posterior becomes

$$P(z \mid \underline{\hat{F}}) \propto \sum_{k} w^{k} c_{1}^{k} \int \tilde{\mathcal{N}}(z, \underline{F} \mid \underline{\eta}^{k\hat{F}}, \mathbf{\Lambda}^{k\hat{F}}) \, d\underline{F}.$$
 (28)

In principle, this integral can be done analytically by moving back to standard parametrisation, i.e., $\Sigma^{k\hat{F}} = \left(\Lambda^{k\hat{F}}\right)^{-1}$ and $\underline{\mu}^{k\hat{F}} = \Sigma^{k\hat{F}}\underline{\eta}^{k\hat{F}}$. The marginalisation can then be done by dropping the corresponding elements from the mean vector and covariance matrix, giving

$$P(z \mid \underline{\hat{F}}) \propto \sum_{k} w^{k} c_{1}^{k} \mathcal{N}(z \mid \underline{\mu}_{z}^{k} \hat{F}, \Sigma_{zz}^{k} \hat{F}).$$
 (29)

Note that this is simply a one-dimensional Gaussian mixture model with a new set of weights given by $w^{k\hat{F}} \equiv w^k c_1^k$.

An important caveat to this result, however, is that the limits of integration are assumed to be $(-\infty, \infty)$; that is, non-physical negative fluxes contribute to the integral. This is the same assumption as used in the derivation in Bovy et al. (2012) using the convolution property of multivariate Gaussians. For this non-blended photo-z, this assumption is sound since the latent fluxes are strongly constrained by the likelihood, meaning that negative fluxes will be strongly down-weighted. However, this will not be the case for the blended photo-z derived in section 4.3 where only the sum of two latent flux vectors is observed.

An alternative approach is to add the boundary prior $\Theta(z, \underline{F})$ as described in section 3.1. This has two effects. Firstly, the prior with this addition must be explicitly normalised, a necessary condition for the model selection. The normalisation factor is given by an integral over the unnormalised prior, i.e.,

$$\mathcal{A}_1 = \iint \Theta(z, \underline{F}) \sum_k w^k \, \mathcal{N}(z, \underline{F} \mid \underline{\mu}^k, \Sigma^k) \, \mathrm{d}z \, \mathrm{d}\underline{F}. \tag{30}$$

This integral can be efficiently estimated using Monte Carlo integration. First, a set of redshifts and fluxes $\{z, \underline{F}\}$ is sampled from the mixture, as described in section 3.4. Since the prior without the boundary prior is normalised to unity as in equation 6, this integral is then equal to fraction of these samples obeying the boundary prior, i.e., where $\Theta(z, \underline{F}) = 1$.

The second effect of adding the boundary prior is that marginalising over fluxes is no longer analytic. Inserting the boundary prior and the corresponding prior normalisation \mathcal{A}_1 , the posterior we want to sample from is given by

$$P(z \mid \underline{\hat{F}}) \propto \mathcal{A}_1 \sum_k w^k c_1^k \int \Theta(z, \underline{F}) \tilde{\mathcal{N}}(z, \underline{F} \mid \underline{\eta}^{k\hat{F}}, \Lambda^{k\hat{F}}) \, d\underline{F}. \quad (31)$$

However, the boundary prior makes this integral nonanalytic and the resulting posterior is not a standard GMM, meaning that it cannot be sampled as described in section 3.4. Instead, we sample from the density given by

$$P(z,\underline{F}\mid\underline{\hat{F}})\propto\mathcal{A}_1\sum_k w^k c_1^k \tilde{N}(z,\underline{F}\mid\underline{\eta}^{k\hat{F}},\pmb{\Lambda}^{k\hat{F}})\,. \tag{32}$$

This is the desired posterior from equation 31 without the marginalisation over fluxes and where we have neglected the

 $^{^{3}}$ This is also referred to as the canonical or information parametrisation.

boundary prior term. This can then be corrected for by rejecting any sample that contains negative fluxes or redshift, leaving only the samples that obey the boundary prior. The marginalisation can then be done trivially by discarding the fluxes and considering only the redshift part of the remaining samples. Since equation 32 is simply a new Gaussian mixture model as before, sampling from this distribution is extremely computationally efficient, as detailed in section 3.4. As described above, the inclusion of the boundary prior is most important for the blended photo-z, though we include it here for completeness and consistency with the blended case later.

4.2 Single-constituent evidence

One of the more computationally demanding aspects of the method of Jones & Heavens (2019) is the use of nested sampling in order to calculate the evidence. A significant advantage of the GMM method presented here is that this expensive integral can be evaluated much more quickly, an important feature for applying the method to future surveys.

The single-constituent evidence \mathcal{E}^1 is defined to be the integral of the unnormalised posterior over the full parameter space, i.e.,

$$\mathcal{E}^{1} = \int \int P(\underline{\hat{F}} \mid \underline{F}) P(z, \underline{F}) d\underline{F} dz.$$
 (33)

As described above, by ignoring the boundary prior, the integral over fluxes can be performed analytically to give a new Gaussian mixture model. Inserting this result into the evidence integral, equation 33 becomes

$$\mathcal{E}^{1} = \sum_{k} w^{k} c_{1}^{k} \int \mathcal{N}(z \mid \underline{\mu}_{z}^{k\hat{F}}, \Sigma_{zz}^{k\hat{F}}) dz.$$
 (34)

Since the multivariate Gaussian density of each component is normalised to unity, the evidence is then given simply by the sum over the new mixture weights, i.e.,

$$\mathcal{E}^1 = \sum_k w^k c_1^k \equiv \sum_k w^{k\hat{F}} \,. \tag{35}$$

In this case, the evidence is analytic and therefore easy to compute. However, as above, computing these integrals analytically implicitly involves contributions from non-physical negative fluxes and redshifts.

To combat this, we can numerically integrate the non-marginalised posterior of fluxes and redshifts including the boundary prior introduced in section 3.1 and the accompanying normalisation from equation 30, i.e.,

$$\mathcal{E}^{1} = \iint \mathcal{A}_{1} \sum_{k} w^{k} c_{1}^{k} \Theta(z, \underline{F}) \tilde{\mathcal{N}}(z, \underline{F} \mid \underline{\eta}^{k\hat{F}}, \boldsymbol{\Lambda}^{k\hat{F}}) \, d\underline{F} \, dz \,. \quad (36)$$

This integral can be evaluated numerically by using fluxes and redshifts sampled from the non-marginalised posterior with the boundary prior removed, given in equation 32. This is another Gaussian mixture model, and thus these samples are computationally efficient to draw, as described in section 3.4. In addition, the posterior samples drawn for inference are also sampled from equation 32 and so can be reused here, saving computation.

Given a set of samples $\{z, \underline{F}\}$ from equation 32, only a fraction \mathcal{F}_1 of these will contain no negative fluxes. Unlike

equation 30, however, this density is not normalised to unity, but rather

$$\begin{aligned} \mathcal{V}_1 &\equiv \mathcal{A}_1 \sum_k w^k c_1^k \iint \tilde{\mathcal{N}}(z, \underline{F} \mid \underline{\eta}^{k\hat{F}}, \mathbf{\Lambda}^{k\hat{F}}) \, \mathrm{d}z \, \mathrm{d}\underline{F} \\ &= \mathcal{A}_1 \sum_k w^k c_1^k \, . \end{aligned} \tag{37}$$

By using this to compute a Monte Carlo estimate of the integral, the evidence can therefore be estimated to be

$$\mathcal{E}^1 \approx \mathcal{V}_1 \mathcal{F}_1 = \mathcal{A}_1 \mathcal{F}_1 \sum_k w^k c_1^k \equiv \mathcal{A}_1 \mathcal{F}_1 \sum_k w^{k\hat{F}} . \tag{38}$$

4.3 Two-constituent posterior

We now extend the inference method to the case of a blended source consisting of two constituent galaxies by deriving the two-constituent posterior. Here, the parameters we wish to infer are the redshifts of each constituent $\{z\} = \{z_1, z_2\}$, given the data vector of observed fluxes \hat{F} .

As before, we start by marginalising over the latent flux vectors. As this is the two-constituent posterior, there are now two flux vectors to marginalise over, $\{\underline{F}\} = \{\underline{F}_1, \underline{F}_2\}$, one for each galaxy. The posterior is therefore given by

$$P(\lbrace z\rbrace \mid \underline{\hat{F}}) = \int P(\lbrace z\rbrace, \lbrace \underline{F}\rbrace \mid \underline{\hat{F}}) \, \mathrm{d}\{\underline{F}\}. \tag{39}$$

Applying Bayes rule, this becomes

$$P(\lbrace z \rbrace \mid \underline{\hat{F}}) \propto \int P(\underline{\hat{F}} \mid \lbrace \underline{F} \rbrace) P(\lbrace z \rbrace, \lbrace \underline{F} \rbrace) d\{\underline{F}\}, \tag{40}$$

where $P(\{z\}, \{\underline{F}\})$ is the joint prior over flux and redshift for both constituents. This prior can be factorised to be written in terms of the individual constituent priors $P(z, \underline{F})$, allowing the GMM to be inserted. However, as described in section 2, the parameters of each constituent are correlated. Thus, the joint prior can be written as

$$P(\lbrace z \rbrace, \lbrace F \rbrace) \propto P(z_1, F_1) P(z_2, F_2) M(z_1, z_2),$$
 (41)

where the blending-related correlations have been factored into a single term

$$M(z_1, z_2) = \pi(z_1, z_2) \left[1 + \xi(z_1, z_2) \right]. \tag{42}$$

Here, $\xi(z_1, z_2)$ is the two-point galaxy correlation function, evaluated at the line-of-sight comoving distance between z_2 and z_1 . This correlation function is commonly modelled as a power law (e.g., Peebles 2001). However, we make no assumption of its form throughout this derivation, requiring only that it can be evaluated given a pair of redshifts. This correlation function was found to have little effect in Jones & Heavens (2019) so the results throughout assume $\xi(z_1, z_2) = 0$. Nevertheless, we include it in the derivations here for completeness. The term $\pi(z_1, z_2)$ represents the sorting condition, given by

$$\pi(z_1, z_2) = \begin{cases} 1 & \text{for } z_1 \le z_2 \\ o & \text{otherwise.} \end{cases}$$
 (43)

The need for these terms is briefly described in section 2; see section 2.3 of Jones & Heavens (2019) for more details.

Any selection effects on the training set are already captured in the prior through the training step. This assumes that the training set is sufficiently representative of the test set, though we note that this caveat applies to machine learning methods in general. The selection effect term of Jones & Heavens (2019) simply acts to disfavour inferring fluxes such that the total flux is near the survey limit, as they are a priori less likely to have been selected. Since the total flux is well constrained by observations, this term has little effect on parameter inferences. Instead, its use is motivated by making the magnitude prior proper. This is necessary for evaluating the marginal likelihood for model comparison. However, our GMM prior is proper by construction. As a result, we do not include the selection effect term here.

As in section 4.1, the model selection requires the joint prior to be normalised. We do this by integrating the prior using Monte Carlo integration. To be able to draw samples from the prior efficiently, we insert the definitions of each term and combine into another Gaussian mixture that can be sampled as described in section 3.4. We also include the boundary prior described in section 3.1 in each constituent prior to prevent contributions to the density from non-physical negative fluxes and redshifts.

Inserting the GMM, correlation and boundary prior terms into equation 41, the joint prior becomes

$$P(\{z\}, \{\underline{\hat{F}}\}) \propto M(z_1, z_2) \Theta(z_1, \underline{F}_1) \Theta(z_2, \underline{F}_2) \sum_k \sum_j w^k w^j \times \\ \mathcal{N}(z_1, \underline{F}_1 \mid \mu^k, \Sigma^k) \mathcal{N}(z_2, \underline{F}_2 \mid \mu^j, \Sigma^j) . \tag{44}$$

We now follow an analogous method to that of section 4.1 by combining the two multivariate Gaussians into a single density. We start by defining a partitioned parameter vector that each density can be written in terms of. This is given by

$$\underline{\theta} = \begin{pmatrix} z_1 \\ \underline{F}_1 \\ z_2 \\ \underline{F}_2 \end{pmatrix} . \tag{45}$$

The product of the densities in equation 44 can then be written as a single Gaussian density in terms of this parameter vector

$$\mathcal{N}(z_1, \underline{F}_1 \mid \underline{\mu}^k, \Sigma^k) \, \mathcal{N}(z_2, \underline{F}_2 \mid \underline{\mu}^j, \Sigma^j) = \mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}), \quad (46)$$

where the new mean vector is given by

$$\underline{\mu}^{kj} = \left(\frac{\underline{\mu}^k}{\underline{\mu}^j}\right) = \begin{pmatrix} \mu_z^k \\ \underline{\mu}_f^k \\ \mu_z^j \\ \mu_f^j \end{pmatrix} \tag{47}$$

and the covariance matrix

$$\Sigma^{kj} = \begin{pmatrix} \Sigma^k & \mathbf{0} \\ \mathbf{0} & \Sigma^j \end{pmatrix} = \begin{pmatrix} \Sigma^k_{zz} & \frac{\Sigma^k}{z_f^k} & 0 & \frac{0}{2} \\ \frac{\Sigma^k}{f_z} & \Sigma^k_{ff} & \frac{0}{2} & \mathbf{0} \\ 0 & \frac{0}{2} & \Sigma^j_{zz} & \Sigma^j_{ff} \\ \frac{0}{2} & \mathbf{0} & \Sigma^j_{fz} & \Sigma^j_{ff} \end{pmatrix}. \tag{48}$$

This combination is trivial since we assume that all correlations between the two constituents have already been factored out into $M(\{z\}, \{\underline{F}\})$. As a result, the two constituent priors are independent and can be combined with the block diagonal covariance matrix defined in equation 48. The joint

prior thus becomes

$$P(\lbrace z \rbrace, \lbrace \underline{\hat{F}} \rbrace) \propto M(z_1, z_2) \Theta(z_1, \underline{F}_1) \Theta(z_2, \underline{F}_2) \times$$

$$\sum_{k} \sum_{j} w^k w^j \mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}),$$

$$(49)$$

i.e., a GMM multiplied by several additional terms. The normalisation of this prior is then given by the integral

$$\mathcal{A}_{2} = \iiint M(z_{1}, z_{2})\Theta(z_{1}, \underline{F}_{1})\Theta(z_{2}, \underline{F}_{2}) \times \sum_{k} \sum_{j} w^{k} w^{j} \mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}) dz_{1} dz_{2} d\underline{\hat{F}}_{\alpha} d\underline{\hat{F}}_{\beta}.$$
(50)

Analogously to equation 30, this can be evaluated using samples drawn from the Gaussian mixture, i.e.,

$$\{z_1,z_2,\underline{F}_1,\underline{F}_2\} \sim G(\underline{\theta}) = \sum_k \sum_j w^k w^j \mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}). \tag{51}$$

Given $n_{\mathcal{A}}$ of these samples $\{z_1^i, z_2^i, \underline{F_1^i}, \underline{F_2^i} \mid i = 1 \dots n_{\mathcal{A}}\}$, we can compute a Monte Carlo integration of \mathcal{A}_2 through importance sampling. Since $G(\underline{\theta})$ is normalised to unity, this integral is given by

$$\mathcal{A}_{2} = \sum_{i} \frac{\left[1 + \xi(z_{1}^{i}, z_{2}^{i})\right] \pi(z_{1}^{i}, z_{2}^{i}) \Theta(z_{1}^{i}, \underline{F}_{1}^{i}) \Theta(z_{2}^{i}, \underline{F}_{2}^{i})}{n_{\mathcal{A}}} \,. \tag{52}$$

If the correlation function and sorting condition were ignored, this would simply be equal to the fraction of samples that obey the boundary prior, as in the definition of \mathcal{A}_1 . Thus, the joint prior is given by

$$P(\lbrace z \rbrace, \lbrace \underline{\hat{F}} \rbrace) = \mathcal{A}_2 M(z_1, z_2) \Theta(z_1, \underline{F}_1) \Theta(z_2, \underline{F}_2) \times$$

$$\sum_k \sum_j w^k w^j \mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}).$$
(53)

This joint prior can then be inserted into equation 40 alongside the definition of the likelihood to develop the posterior. As before, we assume that the likelihood is a multivariate Gaussian centred on the observed fluxes, though we now model the flux as the sum of the constituent fluxes, i.e.,

$$P(\underline{\hat{F}} \mid \underline{F}) = \mathcal{N}(\underline{F}_1 + \underline{F}_2 \mid \underline{\hat{F}}, \Sigma^{\widehat{F}}). \tag{54}$$

Inserting this likelihood and the joint prior into equation 40, the posterior becomes

$$P(\lbrace z \rbrace \mid \underline{\hat{F}}) \propto \mathcal{A}_{2} \iint M(z_{1}, z_{2}) \Theta(z_{1}, \underline{F}_{1}) \Theta(z_{2}, \underline{F}_{2}) \times$$

$$\sum_{k} \sum_{j} w^{k} w^{j} \mathcal{N}(\underline{F}_{1} + \underline{F}_{2} \mid \underline{\hat{F}}, \Sigma^{\hat{F}}) \times$$

$$\mathcal{N}(\underline{\theta} \mid \underline{\mu}^{kj}, \Sigma^{kj}) \, d\underline{F}_{1} \, d\underline{F}_{2}.$$
(55)

To combine the prior term with the likelihood, we rewrite it in terms of natural parameters partitioned in the same way as equation 45. These new parameters are given by

$$\underline{\eta}^{kj} = \begin{pmatrix} \underline{\eta}^k \\ \underline{\eta}^j \end{pmatrix} = \begin{pmatrix} \eta_z^k \\ \underline{\eta}_f^k \\ \eta_z^j \\ \underline{\eta}_f^j \end{pmatrix} \tag{56}$$

12

and

$$\mathbf{\Lambda}^{kj} = \begin{pmatrix} \mathbf{\Lambda}^k & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}^j \end{pmatrix} = \begin{pmatrix} \mathbf{\Lambda}^k_{zz} & \underline{\mathbf{\Lambda}}^k_{zf} & \mathbf{0} & \underline{\mathbf{0}} \\ \underline{\mathbf{\Lambda}}^k_{fz} & \mathbf{\Lambda}^k_{ff} & \underline{\mathbf{0}} & \mathbf{0} \\ \mathbf{0} & \underline{\mathbf{0}} & \mathbf{\Lambda}^j_{zz} & \underline{\mathbf{\Lambda}}^j_{zf} \\ \underline{\mathbf{0}} & \mathbf{0} & \underline{\mathbf{\Lambda}}^j_{fz} & \mathbf{\Lambda}^j_{ff} \end{pmatrix}. \quad (57)$$

By also rewriting the likelihood in terms of the natural parameters $\Lambda^{\hat{F}} \equiv \left(\Sigma^{\hat{F}}\right)^{-1}$ and $\underline{\eta}^{\hat{F}} \equiv \Lambda^{\hat{F}}\underline{\hat{F}}$, the posterior becomes

$$P(\lbrace z \rbrace \mid \underline{\hat{F}}) \propto \iint \mathcal{A}_{2} M(z_{1}, z_{2}) \Theta(z_{1}, \underline{F}_{1}) \Theta(z_{2}, \underline{F}_{2}) \times$$

$$\sum_{k} \sum_{j} w^{k} w^{j} \tilde{N}(\underline{F}_{1} + \underline{F}_{2} \mid \underline{\eta}^{\hat{F}}, \boldsymbol{\Lambda}^{\hat{F}}) \times$$

$$\tilde{N}(\theta \mid \eta^{kj}, \boldsymbol{\Lambda}^{kj}) dF_{1} dF_{2}.$$

$$(58)$$

The two remaining densities can now be combined into a single term given by

$$\tilde{\mathcal{N}}(\underline{F}_1 + \underline{F}_2 \mid \underline{\eta}^{\hat{F}}, \Lambda^{\hat{F}}) \, \tilde{\mathcal{N}}(\underline{\theta} \mid \underline{\eta}^{kj}, \Lambda^{kj}) \propto \tilde{\mathcal{N}}(\underline{\theta} \mid \underline{\eta}^{k\hat{F}}, \Lambda^{k\hat{F}}), \quad (59)$$
 where the combined parameters are given by

$$\underline{\eta}^{k\hat{F}} = \begin{pmatrix} \eta_z^k \\ \eta_f^k + \underline{\eta}^{\hat{F}} \\ \eta_z^j \\ \eta_f^j + \eta^{\hat{F}} \end{pmatrix}$$

$$(60)$$

and

$$\boldsymbol{\Lambda}^{k\hat{F}} = \begin{pmatrix} \boldsymbol{\Lambda}_{zz}^{k} & \boldsymbol{\Lambda}_{zf}^{k} & 0 & \underline{0} \\ \underline{\Lambda}_{fz}^{k} & \boldsymbol{\Lambda}_{ff}^{k} + \boldsymbol{\Lambda}^{\hat{F}} & \underline{0} & \boldsymbol{\Lambda}^{\hat{F}} \\ \underline{0} & \underline{0} & \boldsymbol{\Lambda}_{zz}^{j} & \underline{\Lambda}_{zf}^{j} + \boldsymbol{\Lambda}^{\hat{F}} \end{pmatrix}.$$
(61)

As before, the constant of proportionality c_2^{kj} in equation 59 can also be written in terms of another multivariate Gaussian density

$$c_2^{kj} = \mathcal{N}\left(\mu_{\rm f}^k + \mu_{\rm f}^j \mid \underline{\hat{F}}, \left[\Sigma^{\hat{F}} + \Sigma_{\rm ff}^k + \Sigma_{\rm ff}^j\right]\right). \tag{62}$$

The posterior is thus given by

$$P(\lbrace z\rbrace \mid \underline{\hat{F}}) \propto \iint \mathcal{A}_{2}M(z_{1}, z_{2})\Theta(z_{1}, \underline{F}_{1})\Theta(z_{2}, \underline{F}_{2}) \times \\ \sum_{k} \sum_{j} w^{k} w^{j} c_{2}^{kj} \mathcal{N}(\underline{\theta} \mid \underline{\eta}^{k\hat{F}}, \boldsymbol{\Lambda}^{k\hat{F}}) \, d\underline{F}_{1} \, d\underline{F}_{2}.$$

$$(63)$$

As in the single constituent case, it would be possible to do this integral analytically by ignoring the boundary prior $\Theta(z, \underline{F})$. Converting back to the standard parametrisation, the final posterior would then be given by

$$P(\{z\} \mid \hat{\underline{F}}) \propto \mathcal{A}_2 M(z_1, z_2) \sum_k \sum_j w^k w^j c_2^{kj} \mathcal{N}(z_1, z_2 \mid \underline{\mu}_z^{k\hat{F}}, \Sigma_{zz}^{k\hat{F}}).$$
(64)

With the boundary prior, the integral is no longer analytically tractable. As a result, we take the same approach as in the single constituent case and sample from the full, non-marginalised posterior. An additional complication here are the extra correlations factored into $M(z_1, z_2)$. As a result of this term, the posterior is no longer a Gaussian mixture and therefore does not permit the efficient sampling scheme described in section 3.4.

Instead, we can sample from the full posterior distribution ignoring the contribution of both the the boundary prior and the correlations, modifying the samples post hoc by rejection and reweighting to correct for these respectively. This set of samples is thus drawn from the simplified posterior $H(\theta)$, given by

$$\{z_{1}, z_{2}, \underline{F}_{1}, \underline{F}_{2}\} \sim H(\underline{\theta}) \propto \mathcal{A}_{2} \sum_{k} \sum_{j} w^{k} w^{j} c_{2}^{kj} \times$$

$$\mathcal{N}(\theta \mid \eta^{k\hat{F}}, \mathbf{\Lambda}^{k\hat{F}}).$$

$$(65)$$

This simplified posterior is now a standard GMM, and can therefore be efficiently sampled as described in section 3.4. The neglected terms can now be corrected for separately.

Firstly, the boundary priors can be included by rejecting samples where the flux or the redshift is negative, as in section 4.1. The sorting condition could also be included by simply rejecting samples where it was not respected. However, this is unnecessarily wasteful of computation. Note that mixture component-jk is identical to component-kj under exchange of constituents. Every component is matched with a pair in this way. As a result, the posterior is exactly symmetric, meaning that samples with misordered redshifts can be corrected by simply swapping the order of their constituents.

The redshift correlation function can be corrected for using importance sampling by associating each sample with a weight $[1+\xi(z_1,z_2)]$. All inferences done with these samples would then need to account for these weights. The risk with this importance sampling approach is that regions of parameter space where the correlation function is large could be poorly sampled when using the modified posterior. The effect of the correlation function would then be underrepresented. However, Jones & Heavens (2019) found that including the redshift correlation function when sampling the posterior had little effect on inferences. As a result, we expect any errors from the use of importance sampling here to be negligible.

Given a set of corrected samples of redshift and flux, the marginalisation can then be done in the same way as in section 4.1, by discarding the flux parts of the samples. The distribution of the remaining redshift samples will then be proportional to the marginalised posterior defined in equation 63, as desired.

4.4 Two-constituent evidence

The two-constituent evidence \mathcal{E}^2 is defined as the integral of the blended posterior over both sets of fluxes and redshifts, i.e.

$$\mathcal{E}^2 = \iint P(\underline{\hat{F}} \mid \{\underline{F}\}) \ P(\{z\}, \{\underline{F}\}) \ \mathrm{d}\{z\} \ \mathrm{d}\{\underline{F}\} \ . \tag{66}$$

Inserting the definitions of each term from the full posterior given in equation 63, this expression becomes

$$\mathcal{E}^{2} = \mathcal{A}_{2} \iiint M(z_{1}, z_{2}) \Theta(z_{1}, \underline{F}_{1}) \Theta(z_{2}, \underline{F}_{2}) \times$$

$$\sum_{k} \sum_{j} w^{k} w^{j} c_{2}^{kj} \mathcal{N}(\underline{\theta} \mid \underline{\eta}^{k\hat{F}}, \boldsymbol{\Lambda}^{k\hat{F}}) dz_{1} dz_{2} d\underline{F}_{1} d\underline{F}_{2}.$$

$$(67)$$

As before, we evaluate this integral numerically using Monte Carlo integration. To do this, we can reuse the samples

drawn for the blended posterior inference from $H(\underline{\theta})$ defined in equation 65. Given a set of n_2 of these samples $\{z_1^i, z_2^i, \underline{F}_1^i, \underline{F}_2^i \mid i=1\dots n_2\}$, we can define the weighted fraction

$$\mathcal{F}_{2} = \sum_{i} \frac{[1 + \xi(z_{1}^{i}, z_{2}^{i})] \pi(z_{1}^{i}, z_{2}^{i}) \Theta(z_{1}^{i}, \underline{F}_{1}^{i}) \Theta(z_{2}^{i}, \underline{F}_{2}^{i})}{n_{2}}.$$
 (68)

This is analogous to \mathcal{F}_1 , the fraction of samples drawn from the non-marginalised single-constituent posterior defined in equation 32 that obey the boundary prior, but with the additional blending-related correlations. The simplified posterior $H(\underline{\theta})$ is not normalised to unity. However, the normalisation constant \mathcal{V}_2 is given by the integral over the full support of the distribution, giving

$$\mathcal{V}_{2} \equiv \int \mathcal{A}_{2} \sum_{k} \sum_{j} w^{k} w^{j} c_{2}^{kj} \mathcal{N}(\underline{\theta} \mid \underline{\eta}^{k\hat{F}}, \boldsymbol{\Lambda}^{k\hat{F}}) d\underline{\theta}
= \mathcal{A}_{2} \sum_{k} \sum_{j} w^{k} w^{j} c_{2}^{kj}.$$
(69)

Thus, the two-constituent evidence can be estimated by importance sampling to be

$$\mathcal{E}^2 \approx \mathcal{V}_2 \mathcal{F}_2 = \mathcal{A}_2 \mathcal{F}_2 \sum_k \sum_j w^k w^j c_2^{kj} \,. \tag{70}$$

5 TESTS ON SIMULATED SOURCES

In order to test our method, we construct a two sets of simulated observations to train our model and compare predictions against. These two sets correspond to an LSST-like optical survey (Ivezić et al. 2019), and the same survey with additional Euclid-like infrared observations (Laureijs et al. 2011). The complementarity of LSST and Euclid has been investigated previously (e.g., Rhodes et al. 2017); additional filter bands will help to break colour-redshift degeneracies and therefore enable more accurate photometric redshifts.

Simulated observations are generated by redshifting a template, integrating over the relevant filter response curves, scaling the results to a given i-band magnitude, adding observational noise and imposing selection criteria. We use the set of templates assembled by Coe et al. (2006) containing eight templates. This is the default template set in the commonly used BPZ (Benítez 2000) photometric redshift software.

We randomly sample true redshift, magnitude and template parameters for each source from a prior using emcee (Foreman-Mackey et al. 2013). The single-constituent joint redshift-magnitude-template prior is defined as follows. First we factorise into separate prior terms, i.e.,

$$P(z, m, t) = P(z \mid m)P(t \mid m)P(m), \tag{71}$$

where t is an integer labelling each template and the redshift prior is assumed to be independent of template. The redshift and magnitude priors are then given by the LSST predictions in LSST Science Collaboration et al. (2009). The redshift prior, based on simulated high-redshift galaxy populations (Kitzbichler & White 2007) is given by

$$P(z \mid m) = \frac{1}{2z_0(m)} \left(\frac{z}{z_0(m)}\right)^2 \exp\left(\frac{-z}{z_0(m)}\right),\tag{72}$$

where

$$z_0(m) = 0.0417m - 0.744, (73)$$

and m refers to i-band magnitude. The corresponding i-band magnitude prior, fitted to data from the Canada-France-Hawaii Telescope Legacy Survey (CFHTLS; Hoekstra et al. 2006), is then given by

$$P(m) \propto 10^{0.31(m-25)}$$
 (74)

We also use the template prior from Benítez (2000), given by

$$P(t|m) = f_t \exp\left(-k_t[m - m_0]\right), \tag{75}$$

where we set $m_0 = 20$ and the parameters f_t and k_t , each dependent on the template type, are set to the values given in Benítez (2000).

Once the redshift, magnitude and template are sampled from this joint prior, the intrinsic fluxes are simulated by redshifting the template and integrating over filter response curves. For the optical survey, we use the six LSST filters u, g, r, i, z, Y (LSST Science Collaboration et al. 2009). We use the three Euclid filters Y, J, H (Racca et al. 2016) as additional infrared bands, giving a total of nine bands for the combined surveys.

Finally, we add magnitude-dependent observational noise to each band. For the optical bands, this is given by the predicted LSST noise model (LSST Science Collaboration et al. 2009). The 5σ depth of point sources in the Euclid Y, J and H bands is 24mag (Laureijs et al. 2011), the same depth as point sources in the LSST i-band (LSST Science Collaboration et al. 2009). We therefore approximate the observational noise in the Y, J and H bands by assuming that their signal-to-noise is equal to that of the i-band.

In order to simulate the flux of blended sources, we add the intrinsic fluxes of two simulated sources and add observational noise corresponding to the total blended flux. The two-constituent prior also needs to account for the blended-related terms described above. The redshift prior includes the sorting condition $\pi(z_1, z_2)$, though we assume no clustering, i.e., $\xi(z_1, z_2) = 0$, as it has a negligible effect at large separations when $z_1 \not\approx z_2$. We also impose a prior on the faintest *i*-band magnitude of either constituent such that it must be brighter than a 5σ detection. A cut like this is necessary since it only makes sense to consider a source blended when each constituent is sufficiently bright. If a constituent is too faint, it should instead be considered to be a contributor to the background flux, rather than that of the source itself.

Finally, we select sources by imposing an *i*-band magnitude cut of $m_i < 25$. This corresponds to the LSST gold sample (LSST Science Collaboration et al. 2009), a population of $\approx 4 \times 10^9$ high signal-to-noise galaxies. For each of the two sets of simulated sources, we randomly select 10000 single-constituent sources to act as a training set, a further 10000 single-constituent sources for the unblended test set, and 10000 two-constituent sources for the blended test set.

Given the unblended training set, we use the procedure described in section 3.3 to set the number of mixture components N. Using 3-fold cross-validation, we test from N=5 to N=100 in multiples of 5, measuring the RMS scatter $\sigma_{\rm RMS}$ defined in equation 14 at each iteration. In order to

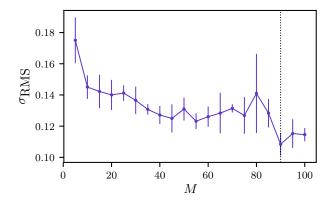


Figure 3. Results of the cross-validation for the LSST-like simulated data. The points show the RMS scatter averaged over the three folds, while the error bars show the error on the mean. We choose the number of components to be N=90, minimising the average RMS scatter as indicated by the dotted black line.

evaluate this, we must define a way to calculate a point estimate $z_{\rm p}$ from a set of n_2 samples $\{z_{{\rm p},i} \mid i=1\dots n_2\}$ drawn from the posterior defined in section 4.1. We therefore define this point estimate to be the mean of these samples, as this is equivalent to a Monte Carlo estimate of the expectation value of the redshift, i.e.,

$$z_{\rm p} \equiv \frac{1}{n_2} \sum_{i=1} z_{{\rm p},i} \approx \int P(z \mid \hat{\underline{F}}) z \, dz \,.$$
 (76)

The results of this cross-validation are shown in Fig. 3. We find the average RMS scatter across all folds $\overline{\sigma_{\rm RMS}}$ to be minimised when N=90 with $\overline{\sigma_{\rm RMS}}=0.108$. We therefore use a mixture comprised of 90 components to fit the entire training set for use throughout.

Examples of one-constituent posteriors inferred using samples from the distribution defined in section 4.1 and conditioned on the LSST-like data are shown in Fig. 4. The four panels in this figure show the variety of shapes of posteriors that can result from photometric redshifts and can be represented by the GMMs presented here.

The top two panels of Fig. 4 shows examples of well constrained, accurate posteriors; their shapes are symmetric and close to that of a single Gaussian. However, the posterior shown in the bottom left panel is left-skewed. This long-tailed posterior is a common occurrence in the results of photometric redshift inference. Despite being very non-Gaussian, it can be represented by a mixture of components. Finally, the bottom right panel shows an example of a bimodal posterior that can be easily represented by a mixture of well separated components. While the true redshift is contained well within the lower peak of this posterior, the bimodality has pulled the mean redshift to between the two peaks. As a result, the point estimate is inaccurate, despite the true redshift lying at a point of significantly non-zero posterior density. This demonstrates the loss of information resulting from the compression of a full posterior distribution to a single point estimate.

Examples of two-constituent posteriors inferred using samples from the distribution defined in section 4.3 are

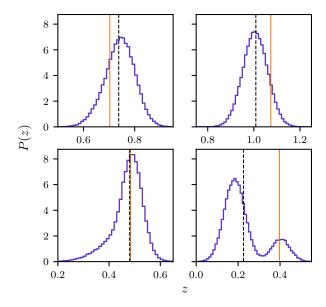


Figure 4. Plot showing four examples of single-constituent posteriors sampled using our method on the unblended LSST-like data. The black dashed lines indicate the sample means we use to define the point estimates $z_{\rm p}$. The true redshifts are indicated by the solid orange lines.

shown in Fig. 5. These samples are also drawn from posteriors conditioned on the LSST-like data.

The left panel of Fig. 5 shows a well constrained posterior. One edge of the joint distribution lies along the $z_1 = z_2$ line. As a result, the effect of the sorting condition $\pi(z_1, z_2)$ can be seen clearly, sharply cutting the joint distribution. The centre panel shows a joint posterior that results in highly skewed marginal distributions. As before, the long tail of the z_2 marginal distribution pulls the mean redshift away from the peak. This demonstrates that, since point estimates are inevitably less informative than the full posterior distribution, the choice of how these point estimates are defined can significantly alter their accuracy. In this case, the accuracy of the point estimate would be increased by choosing z_2 to be the redshift where the posterior peaks, i.e., the maximum a posteriori (MAP) value. However, we found that MAP point estimates were less accurate over the whole sample on average. Finally, the right panel of Fig. 5 shows an example of a highly multimodal posterior that can arise in the two-constituent case.

While less informative than the full posterior distributions, point estimates are still a common product of photometric redshift inference. A plot of these point estimates, defined as the mean of samples drawn from the posterior, against the true redshift for single-constituent data from the two simulated surveys is shown in Fig. 6.

This figure shows that the method performs well in the single-constituent case, i.e., on the standard photometric redshift inference problem. The vast majority of sources have their redshifts recovered accurately; this can be seen by the significant density of points around the $z_{\rm p}=\hat{z}_{\rm s}$ line, demonstrated in the plot by the colour of the points. Comparing the panels for the two simulations, the most significant dif-

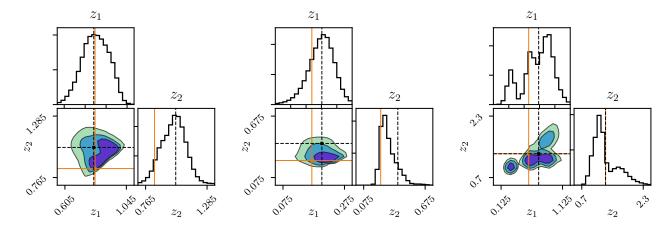


Figure 5. Plot showing three examples of two-constituent posteriors sampled using the GMM on the blended LSST-like data. The black dashed lines indicate the sample means we use to define the point estimates z_p . The true redshifts of each constituent are indicated by the orange lines.

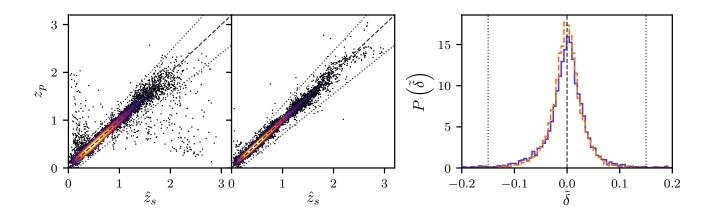


Figure 6. Plot showing the point-estimate results obtained from the GMM on the unblended simulated data. The left and right scatter plots show the point estimate results for the LSST-like and the combined LSST-Euclid-like surveys respectively. These plots show the benefit of additional bands and increased wavelength coverage from near-infrared data in reducing outliers. The dashed line denotes $z_p = \hat{z}_s$, and the dotted lines indicate our outlier definition where $|z_p - \hat{z}_s| \ge 0.15(1 + \hat{z}_s)$. Points are coloured according to their density on the scatter plots to illustrate overplotting. The right panel shows the distribution of the normalised error $\tilde{\delta}$, defined in equation 13. The solid purple line shows the results for the LSST-like survey, while the orange dashed line shows the results for the combined LSST-Euclid-like survey. The black dashed and dotted lines are defined as in the scatter plots.

ference is in the number of outliers, which is reduced in the simulations with additional infrared data. This can also be seen in the third panel, a histogram of the reduced error $\tilde{\delta}$ defined in equation 13. When zoomed around the majority of values at small errors, the difference between the histograms for the two sets of simulations is negligible.

This reduction of outliers is expected, as the additional filters can help to lift the colour-redshift degeneracies discussed in section 3. We define outliers to be sources where $|z_{\rm p}-\hat{z}_{\rm s}|\geq 0.15(1+\hat{z}_{\rm s})$. This boundary is shown as a dotted line in Fig. 6.

In order to quantify the accuracy of these point estimates, we can use several metrics. Firstly, we use the RMS scatter defined in equation 14. We find this scatter to be $\sigma_{\rm RMS} = 0.105$ for the LSST-like simulations, and $\sigma_{\rm RMS} = 0.038$ for the simulations with additional infrared

data. While this difference is significant, it is primarily driven by the reduction of outliers by the infrared data.

In the LSST-like survey, 1.82% of sources are outliers. This is reduced to 0.10% in the combined LSST-Euclid-like simulations. These outliers have significant errors by definition, are therefore can have a significant effect on the measured RMS scatter. In order to identify these outliers as the most significant driver of the difference in accuracy between the two sets of simulations, we measure the RMS scatter while neglecting these sources, as in the photometric redshift accuracy tests of Hildebrandt et al. (2010). When this is done, the RMS of the LSST-like simulations drops to $\sigma_{\rm RMS}=0.036$, while the scatter of the simulations with additional Euclid-like data becomes $\sigma_{\rm RMS}=0.031$. Since these values are now far closer and the latter change was less dra-

matic, we conclude that the biggest benefit afforded by the additional bands is the reduction of outliers.

We also evaluate the same metrics on point estimates of the redshifts of the blended simulated data. These point estimates are defined to be the mean of posterior samples, as in the single-constituent case. A plot of these point estimates for each set of simulated data is shown in Fig. 7.

The blended redshift inference is a more challenging problem than standard photometric redshifts of unblended sources. However, while the scatter plots in Fig. 7 are noisier than the single-constituent plots in Fig. 6, many redshifts are still recovered accurately. This can be seen in the high density of points around $z_{\rm p}=\hat{z}_{\rm s}$, again demonstrated by their colour. This increase in noise over the single-constituent case is expected, as the same number of data-points per source are used here to constrain twice the number of parameters.

As in the single-constituent case, the addition of additional bands in the infrared reduces both the RMS scatter and the number of outliers. For the LSST-like survey, we find the scatter to be $\sigma_{\rm RMS}=0.171$, while the combined LSST-Euclid-like survey has a scatter of $\sigma_{\rm RMS}=0.145$. The outlier rate of the former survey is found to be 17.5%, while that of the latter is reduced to 12.4%.

As discussed in section 3, an important part of the results of photometric redshift inference are PDFs. Unlike simple point estimates, PDFs represent the full statistical knowledge of the redshift being inferred and are essential for rigorously propagating uncertainties. It is therefore also important that the quality of the resulting PDFs are assessed.

A conceptual problem with assessing the quality of PDFs is that there is no true PDF that they can be compared against. This is in contrast to point estimates where the spectroscopic redshift provides a known ground truth against which to compare. Instead, Wittman et al. (2016) introduce a frequentist method to test the widths of PDFs that relies on credible intervals (CIs).

The definition of CIs follows directly from that of posterior PDFs. For a given posterior $P(\theta \mid d)$ that is correctly normalised, the conditional probability that the parameter θ will lie within an interval $[\theta_{\text{low}}, \theta_{\text{high}}]$ is given by the integral of the posterior over that interval, i.e.,

$$P(\theta_{\text{low}} \le \theta \le \theta_{\text{high}} \mid d) = \int_{\theta_{\text{low}}}^{\theta_{\text{high}}} P(\theta \mid d) \, d\theta \,. \tag{77}$$

The CI corresponding to a particular percentage is then defined to be the interval over which equation 77 equals this percentage. In general, this interval will not be unique, since the integral over many different intervals can be the same. For this reason, the credible interval is often defined to be the highest posterior density (HPD) interval, the interval covering the shortest length in parameter space for a given integral. In general, this region does not need to be contiguous; the HPD region of multimodal posteriors will instead be made up of several subintervals.

A conceptually simple way to define this HPD region is to consider a horizontal line spanning the entirety of parameter space, drawn on a plot of the PDF. As this line is moved downwards, it will begin to intersect the PDF. The regions between these intersections can then be integrated to give an area. The intervals contained within these intersections are the HPD region corresponding to this area. Since this area will monotonically increase as the line is moved downwards, this provides a way to define the HPD region for a given percentage CI.

An intuitive interpretation of these intervals is that, given many repetitions of the experiment and the subsequent construction of many such intervals of area α , the true parameter would be contained within a fraction α of these intervals. This notion is the interpretation of frequentist confidence intervals as coverage probabilities. However, while this interpretation is intuitive, it is not guaranteed by a Bayesian analysis. Instead, posteriors where this coverage probability property holds are said to be *calibrated*, and several methods having been proposed to calibrate posteriors (e.g. Syring & Martin 2018; Sellentin & Starck 2019).

The method introduced in Wittman et al. (2016) tests whether the posteriors resulting from a photometric redshift method are calibrated. If they are, we should expect that 50% of sources have their true redshift within their 50% CI. The equivalent statement can be made for all levels of CI, generalising this to a continuous test. The method may therefore give an indication of the performance of the method, and such a test has been widely adopted in the photometric redshift literature (e.g., Leistedt & Hogg 2017; Gomes et al. 2017; Duncan et al. 2018; Meshcheryakov et al. 2018; Amaro et al. 2018; Rodríguez-Muñoz et al. 2019).

By definition, if the true redshift of a source lies within its 50% CI, it will also lie within all CIs corresponding to larger percentages, as the 50% CI will be a subset of these. It is therefore sufficient to measure only the threshold CI that just contains the true redshift. This will have one of the interval edges at the true redshift. This region can therefore be measured by drawing the horizontal line detailed above so that it intersects the posterior at the true redshift. The area c corresponding to this interval is measured for each galaxy in the sample being tested. The cumulative distribution function (CDF) of these areas CDF(c) can then be calculated. Wittman et al. (2016) note that for calibrated posteriors, the plot of this CDF against areas should be diagonal, i.e., CDF(c) = c. The deviation away from this line therefore measures how overconfident or underconfident the PDFs are.

A plot of this test for the LSST-like simulated data is shown in Fig. 8. This figure shows that both the one-and two-constituent posteriors are approximately calibrated and their CIs can therefore be interpreted in a frequentist manner.

Finally, Fig. 9 shows the relative probability for the blended and unblended models $\mathcal{P}_{2,1}$ calculated for the blended data of both simulated surveys. This quantity is calculated using the evidences derived in sections 4.2 and 4.4 using equation 3. We assume a ratio of model priors of unity, i.e., we do not a priori favour either the one- or two-constituent models. A blended source is then favoured when $\ln \mathcal{P}_{2,1} > 1$. We find that the LSST-like survey identifies 92.4% of blended sources, while the survey with additional infrared data identifies 89.3%.

6 GAMA BLENDED SOURCES CATALOGUE

In addition to the simulated observations presented in section 5, we also test our method against real observations.

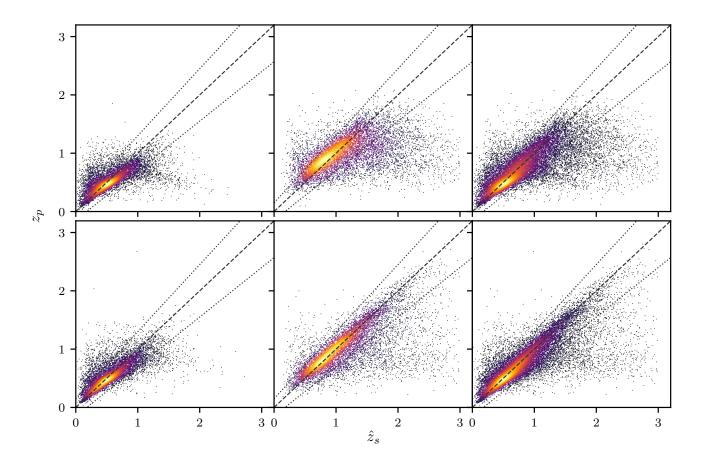


Figure 7. Plot showing the point-estimate results obtained from the GMM on the blended simulated data. The top row shows the results for the LSST-like survey, and the bottom row shows results for the combined LSST-Euclid-like survey. The left plots show $z_{\rm p,1}$, the point estimate of the redshift for the lower-redshift constituent in each blended source. The centre plots show $z_{\rm p,2}$, corresponding to the higher-redshift constituent in each blended source. The right plots combine both $z_{\rm p,1}$ and $z_{\rm p,2}$. The dashed lines denotes $z_{\rm p} = \hat{z}_{\rm s}$, and the dotted lines indicate our outlier definition where $|z_{\rm p} - \hat{z}_{\rm s}| \ge 0.15(1 + \hat{z}_{\rm s})$. Points are coloured according to their density on the scatter plots to illustrate overplotting.

To do this, we use data from the Galaxy And Mass Assembly (GAMA) survey (Baldry et al. 2017), a spectroscopic survey of > 150 000 sources. Alongside this spectroscopy, these sources were also imaged in optical wavelengths by the Sloan Digital Sky Survey (SDSS) (Stoughton et al. 2002) and in infrared wavelengths by the VISTA Kilo-degree Infrared Galaxy (VIKING) Survey (Edge et al. 2013). Hill et al. (2011) used this imaging data to create self-consistent, aperture-matched photometry in nine bands u, g, r, i, z, Y, J, H, K for all sources within the GAMA survey. As a result, these sources have both high-quality photometry and accurate spectroscopic redshifts for training and testing our photometric redshift method.

Holwerda et al. (2015) used this data to spectroscopically identify blended sources in order to search for stronglens candidates. The resulting GAMA blended sources catalogue contains blended photometry for 280 sources, alongside the spectroscopic redshift of each constituent. We therefore use this catalogue to test the performance of our method on real observations of blended sources. To accompany this, we also randomly select two sets of 10000 unblended sources for a training and test set.

As for the simulated observations, we use 3-fold cross-validation to find the number of mixture components N that minimises $\overline{\sigma_{\rm RMS}}$ the RMS scatter averaged over all folds. The results of this are shown in Fig. 10. We find the minimum scatter when the number of mixture components is N=45, giving $\overline{\sigma_{\rm RMS}}=0.066$. We therefore continue with a GMM of 45 components fitted to the 10000 unblended training sources.

We then compute point estimates of the single-constituent redshifts by averaging samples drawn from the posterior as before. A plot of this is shown in Fig. 11. We find the RMS scatter to be $\sigma_{\rm RMS}=0.067$, with 3.6% of sources being outliers.

A scatter plot of the two-constituent point estimates is shown in Fig. 12. As in the simulated case, the blended results are noisier than the single-constituent case. We find the RMS scatter to be $\sigma_{\rm RMS}=0.091$, and 10.8% of sources to be outliers.

Examples of single-constituent posteriors are shown in Fig. 13. Like the single-constituent posteriors conditioned on the simulated data, these distributions show a variety of shapes. However, the posteriors for the GAMA data are significant.

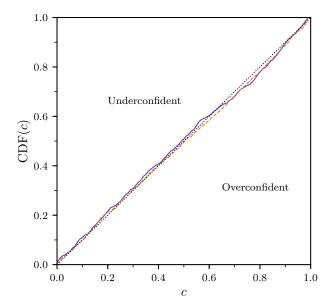


Figure 8. Plot showing the results of the posterior width test performed on posteriors obtained from our method on LSST-like simulated data. The solid purple line shows the results for the single-constituent posteriors, and the dashed orange line shows the results for the two-constituent posteriors. The black dotted line indicates the result where posteriors are calibrated, while lines that go above and below this indicate posteriors that are wider and narrower than calibrated posteriors respectively.

nificantly less multimodal. This is likely because the GAMA sources are, on average, lower redshift than the simulated sources. The main cause of the bimodality in the simulated case is the colour-redshift degeneracy described in section 3, which low- and high-redshift sources to be confused. However, high redshifts are a priori very unlikely here, as they do not appear in the training set. As a result, these higher redshift peaks are significantly disfavoured.

The same lack of multimodality is also exhibited in the blended posteriors conditioned on the GAMA data. Examples of these are shown in Fig. 14. These posteriors show a variety of non-Gaussian shapes as in the simulated case, with many of the marginal redshift distributions displaying long tails. The joint distribution in the left panel of Fig. 14 also shows the hard cut resulting from the sorting condition $\pi(z_1, z_2)$, as the left panel of Fig. 5 does.

Fig. 15 shows the plot testing the posterior widths for both the one- and two-constituent posteriors. As in the simulated case, the one-constituent posteriors are very close to being calibrated. However, the CDF for the two-constituent posteriors lies significantly below the diagonal, suggesting that the posteriors are overconfident, i.e., they are too narrow. As discussed above, while it is not guaranteed that Bayesian CIs provide frequentist coverage probabilities, this suggests that there are features on the flux-redshift relation of the blended constituents that are not captured by the model trained on the unblended training data.

This interpretation is supported by Fig. 16 which shows the inferred relative probability of sources from the blended sources catalogue being blended and unblended $\mathcal{P}_{2,1}$. Here,

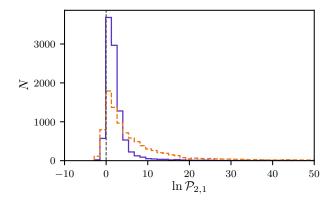


Figure 9. Histograms of the log of the relative probabilities for the blended and unblended models obtained using Bayesian model comparison on the simulated blended data. The solid purple histogram shows the result for the LSST-like survey, while the dashed orange histogram shows the result for the combined LSST-Euclid-like survey. The black dashed line indicates no preference for either the unblended or blended model. Larger values of $\mathcal{P}_{2,1}$ favour the blended model more.

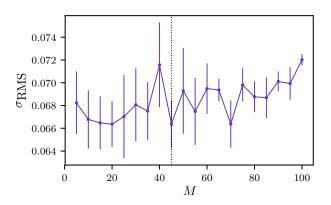


Figure 10. Results of the cross-validation for the GAMA blended sources catalogue data. The points show the RMS scatter averaged over the three folds, while the error bars show the error on the mean. We choose the number of components to be N=45, minimising the average RMS scatter as indicated by the dotted black line.

only 33.4% of blended sources are correctly identified as blended by having $\mathcal{P}_{2,1} > 1$. While the redshifts are reasonably well-recovered, the Bayesian model selection will disfavour a more complicated model when the improvement in the fit is insufficient. As above, this suggests a difference between the blended and unblended constituents.

We can test for a difference between the blended and unblended constituents by incrementally removing sources where this difference is greatest and checking whether this leads to an improvement in the summary statistics. We therefore require a quantity to probe the representativeness of a given vector of fluxes. For this, we consider the density

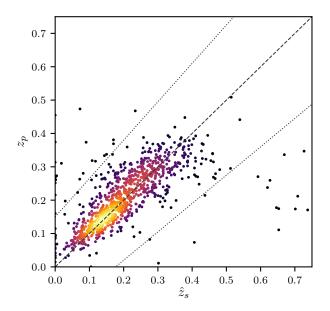


Figure 11. Plot showing the point-estimate results obtained from the GMM on the unblended GAMA data. The dashed line denotes $z_{\rm p}=\hat{z}_{\rm s}$, and the dotted lines indicate our outlier definition where $|z_{\rm p}-\hat{z}_{\rm s}| \geq 0.15(1+\hat{z}_{\rm s})$. Points are coloured according to their density on the scatter plots to illustrate overplotting.

ratio

$$\mathcal{R}(\underline{F}) = \frac{P_{\text{test}}(\underline{F})}{P_{\text{train}}(\underline{F})},\tag{78}$$

where $P_{\text{train}}(\underline{F})$ is the density of fluxes is the training set, $P_{\text{train}}(\underline{F})$ is the density of fluxes is the test set and \underline{F} is the flux vector at which both of these densities are evaluated.

In order to estimate this ratio, we use the nearest-neighbour method of Kremer et al. (2015). The method first considers the training set, and measures the hypervolume that contains the $n_{\rm nei}$ nearest neighbours of a flux \underline{F} . The number of test-set samples $n_{\rm test}(\underline{F})$ within that hypervolume centred on \underline{F} is then counted. The estimate for the density ratio is then given as the ratio of these counts, i.e.,

$$\mathcal{R}(\underline{F}) \approx \frac{n_{\text{nei}}}{n_{\text{test}}(F)}$$
 (79)

This nearest-neighbour method for estimating the density ratio was first presented in Lima et al. (2008), and was used to estimate the redshift distribution of a photometric galaxy sample by weighting spectroscopic galaxies. However, the accuracy of this method depends on $n_{\rm nei}$, the number of neighbours considered. If $n_{\rm nei}$ is too large, the density ratio is estimated over too large a volume, while an estimate where $n_{\rm nei}$ is too small will be dominated by statistical errors. To this end, Kremer et al. (2015) present a model-selection method based on cross-validation to optimise $n_{\rm nei}$.

As discussed throughout this paper, a complication of blended sources is that the flux of each constituent is not observed independently, only the blended combination. As a result, the destiny ratio must be evaluated using constituent fluxes sampled from the marginal posterior $P(\underline{F}_n|\hat{\underline{F}})$, where \underline{F}_n is the flux of constituent n. As described in section 4.3, this can be accomplished by sampling from the simplified

posterior defined in equation 65, and rejecting samples that do not obey the boundary prior. The marginalisation over all redshifts and the flux of the other constituent can than be done by simply ignoring these elements of the sampled vectors.

Given a set of $n_{\rm F}$ flux samples $\{\underline{F}_n^i \mid i=i\dots n_{\rm F}\}$ from constituent n, we evaluate the density ratio $\mathcal{R}(\underline{F})$ for each sample and average the result to give the expectation value

$$E[\mathcal{R}(\underline{F})] \equiv \int \mathcal{R}(\underline{F}_n) P(\underline{F}_n | \hat{\underline{F}}) \, d\underline{F}_n \approx \frac{1}{n_F} \sum_i \mathcal{R}(\underline{F}_n^i). \tag{80}$$

This expectation value is the quantity we use to estimate the representativeness of blended constituents. This allows us to test for differences between the blended and unblended constituents. To do this, we keep sources in our sample only if the expectation of the density ratio for both of their constituents is over a threshold value $\mathcal{R}_{\rm th}$, i.e., sources that obey

$$\frac{\mathbb{E}[\mathcal{R}(\underline{F}_n^i)]}{\max(\mathbb{E}[\mathcal{R}(\underline{F})])} \ge \mathcal{R}_{\text{th}}, \quad n \in \{1, 2\},$$
(81)

where we have normalised the expectation values by $\max(E[R(F)])$, the maximum expectation value over both constituents of all sources.

Fig. 17 shows the change in summary statistics as the threshold ratio is increased. As expected, the RMS scatter and number of outliers are both reduced as this ratio is increased, at the expense of more sources being removed from the sample. This effect can also be seen in the lower two rows of Fig. 12, where the effects of two different threshold values on the point estimates are compared with the unmodified results. When the threshold is set at $\mathcal{R}_{\rm th} = 0.45$ as in the centre row, the RMS scatter has been reduced to $\sigma_{\rm RMS} = 0.078$, while the percentage of sources that are outliers has reduced to 5.97%. At this level, 70.7% of sources remain in the sample. By increasing the threshold to $\mathcal{R}_{\mathrm{th}} = 0.8$ as in the bottom row, the RMS scatter and percentage of outliers decrease to $\sigma_{\rm RMS} = 0.077$ and 4.34% respectively. These are modest improvements over the less strict threshold, but come at the cost of leaving only 40.9% of sources remaining in the sample.

These results demonstrate the importance of representative training sets. Differences between the training and test sets, often referred to as covariate shift, are a general problem for machine learning-based methods that obtain all of their information from the training set. A possible cause of differences here is that surveys select sources based on a magnitude cut, imparting selection effects on the sample. Since blended sources will be selected based on their total blended flux, blended constituents can be fainter than those that are unblended. The simulated sources presented in section 5 are selected in this way and so contain this effect. However, the intrinsic properties of galaxies vary with magnitude, meaning that the test set could contain faint constituents that have no corresponding examples in the test set. Selection effects imparted by the selection criteria of sources in the blended sources catalogue, such as certain redshift differences being easier to select spectroscopically, are also not accounted for here.

One solution to this problem is to improve the training set so that it is more representative. By including sources in the training set fainter than the magnitude limit of the

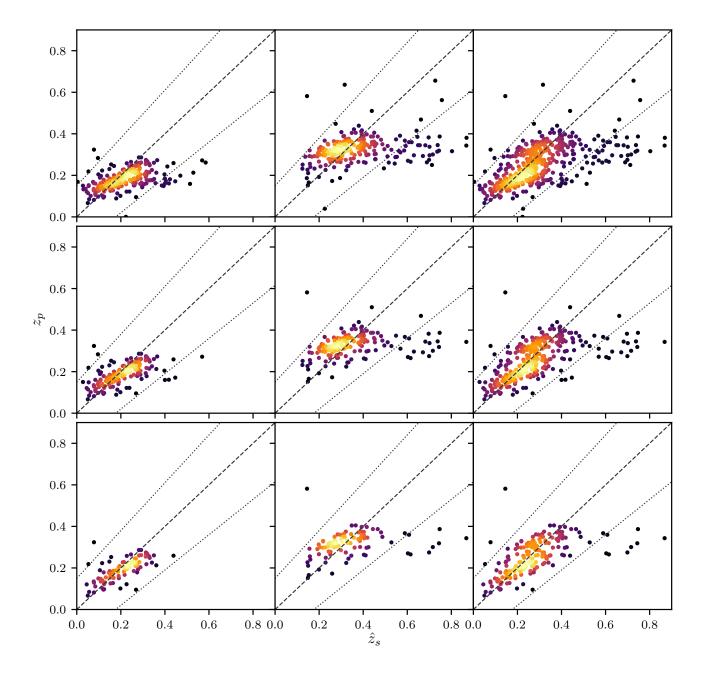


Figure 12. Plot showing the point-estimate results obtained from the GMM on the data from the GAMA blended sources catalogue, with various density ratio thresholds. The left column shows $z_{\rm p,1}$, the point estimate of the redshift for the lower-redshift constituent in each blended source. The centre column shows $z_{\rm p,2}$, corresponding to the higher-redshift constituent in each blended source. The right column combines both $z_{\rm p,1}$ and $z_{\rm p,2}$. The top row shows the results for the full sample, while the centre and bottom rows have sources with expected density ratios less than 0.45 and 0.8 removed respectively. where the expected density ratio is defined in equations 80 and 81. Imposing this density ratio threshold removes sources that are least well-represented in the training set, and so we would expect the results to improve as the threshold is increased. As indicated in the text, the summary statistics improve as expected by making these cuts. This can also be seen visually in this figure by comparing the lower two rows with the full sample in the top row. The dashed lines denotes $z_{\rm p} = \hat{z}_{\rm s}$, and the dotted lines indicate our outlier definition where $|z_{\rm p} - \hat{z}_{\rm s}| \ge 0.15(1 + \hat{z}_{\rm s})$. Points are coloured according to their density on the scatter plots to illustrate overplotting.

test set, the model can learn the faint-end flux-redshift relation. The selection effects of blended sources could also be learned directly by training using a blended training set as described in section 3.1. However, as detailed above, assembling a representative blended training set in practice

could be difficult. For the tests presented here, the GAMA blended sources catalogue contains far too few sources to be amenable to fitting in this way.

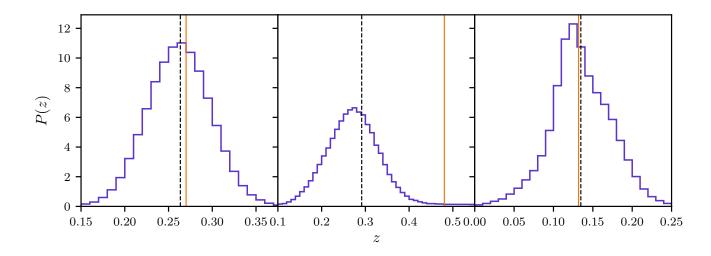


Figure 13. Plot showing three examples of single-constituent posteriors sampled using the GMM on the unblended GAMA data. The black dashed lines indicate the sample means we use to define the point estimates z_p . The true redshifts are indicated by the orange lines

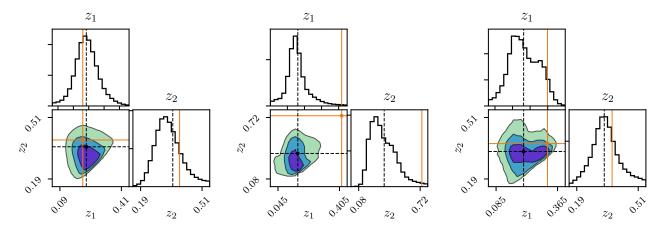


Figure 14. Plot showing three examples of two-constituent posteriors sampled using the GMM on data from the GAMA blended sources catalogue. The black dashed lines indicate the sample means we use to define the point estimates z_p . The true redshifts of each constituent are indicated by the orange lines.

7 CONCLUSIONS

Future galaxy surveys will observe to unprecedented depths in order to drive their increases in precision of cosmological constraints. However, these improvements to constraints on cosmological parameters will be accompanied by several new complications to the analysis. The increased number density of sources will increase both the number of sources that are blended and the total number of sources observed.

This paper presents a photometric redshift method for blended sources based on Gaussian mixture models. Using these models, our method learns the flux-redshift distribution from a set of unblended training galaxies. This choice of model permits the derivation of posteriors that can be sampled efficiently, allowing the method to scale to large samples. By using Bayesian model selection techniques, this method can also infer the number of constituents within a blended sources efficiently.

This work extends previous uses of GMMs in photometric redshift applications Bovy et al. (2012) to the case of blended sources. It also extends the template-based method to infer the redshifts of blended sources directly from their blended photometry first introduced in Jones & Heavens (2019). The method described therein relies on nested sampling for inference and so will not scale to the large sample sizes of future galaxy surveys such as LSST (Ivezić et al. 2019). The method presented in this paper is significantly faster, making it suitable for these upcoming surveys. Many modern methods of photometric redshifts are machine learning-based, as training these methods on a representative training set can allow them to achieve very high accuracy and avoid the problems associated with small template sets. This paper extends the blended photometric redshift method of Jones & Heavens (2019) to this data-driven approach.

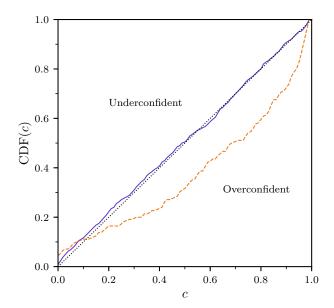


Figure 15. Plot showing the results of the posterior width test performed on posteriors obtained from our method on GAMA data. The solid purple line shows the results for the single-constituent posteriors, and the dashed orange line shows the results for the two-constituent posteriors. The black dotted line indicates the result where posteriors are calibrated, while lines that go above and below this indicate posteriors that are wider and narrower than calibrated posteriors respectively.

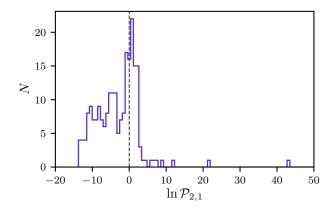


Figure 16. Histogram of the log of the relative probabilities for the blended and unblended models obtained using Bayesian model comparison on the blended GAMA data. The black dashed line indicates no preference for either the unblended or blended model. Larger values of $\mathcal{P}_{2,1}$ favour the blended model more.

The accuracy of all machine learning-based photometric redshift methods is dependent of the training set. Using training sets that are unrepresentative could result in redshift inferences that are biased and posterior distributions that are too narrow. In cases where unblended galaxies are not representative of individual components in a blended source, potentially as a result of selection effects, our method can generalise to learn the blended flux-redshift relation di-

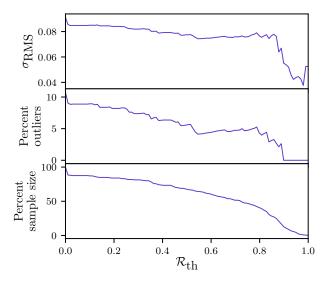


Figure 17. Plot showing the change in summary statistics for the GAMA blended sources as the density ratio threshold $\mathcal{R}_{\rm th}$ is increased. The top panel shows the RMS scatter $\sigma_{\rm RMS}$. The centre panel shows the percentage of sources that are outliers, defined as $|z_{\rm p}-\hat{z}_{\rm s}|\geq 0.15(1+\hat{z}_{\rm s})$. The bottom panel shows the percentage of sources remaining from the original sample after the threshold has been applied.

rectly from blended training data. While this naturally accounts for differences between blended and unblended galaxies, it also increases the size of the required training set.

The method presented here represents a different approach to analysing blended sources than is currently used. Rather than separating blended observations into separate constituents, we infer the redshifts jointly for all constituents. As a result, our method naturally captures uncertainties and correlations which can be difficult to estimate for deblending-based analyses. This approach could be extended to other quantities of interest for cosmological analysis such as galaxy shapes by constructing forward models of source images. By doing this, correlations associated with blending can be propagated fully throughout the rest of the analysis, providing the best understanding of uncertainties on cosmological constraints.

ACKNOWLEDGEMENTS

We thank Andrew Jaffe, Daniel Mortlock and Boris Leistedt for helpful discussions, and the anonymous referee for many useful suggestions that have improved this paper. DMJ acknowledges funding from STFC through training grant ST/N504336/1. GAMA is a joint European-Australasian project based around a spectroscopic campaign using the Anglo-Australian Telescope. The GAMA input catalogue is based on data taken from the Sloan Digital Sky Survey and the UKIRT Infrared Deep Sky Survey. Complementary imaging of the GAMA regions is being obtained by a number of independent survey programmes including GALEX MIS, VST KiDS, VISTA VIKING, WISE, Herschel-ATLAS, GMRT and ASKAP providing UV to radio coverage. GAMA

is funded by the STFC (UK), the ARC (Australia), the AAO, and the participating institutions. The GAMA website is http://www.gama-survey.org/. Based on observations made with ESO Telescopes at the La Silla Paranal Observatory under programme ID 179.A-2004.

REFERENCES

```
Abbott T. M. C., et al., 2018, MNRAS, 480, 3879
```

Almosallam I. A., Jarvis M. J., Roberts S. J., 2016, MNRAS, 462, 726

Amaro V., et al., 2018, Monthly Notices of the Royal Astronomical Society, 482, 3116

Baldry I. K., et al., 2017, preprint, (arXiv:1711.09139)

Benítez N., 2000, ApJ, 536, 571

Benítez N., et al., 2009, ApJ, 692, L5

Bolzonella M., Miralles J.-M., Pelló R., 2000, A&A, 363, 476

Bovy J., Hogg D. W., Roweis S. T., 2011, Annals of Applied Statistics, 5, 1657

Bovy J., et al., 2012, ApJ, 749, 41

Carliles S., Budavári T., Heinis S., Priebe C., Szalay A. S., 2010, ApJ, 712, 511

Carrasco Kind M., Brunner R. J., 2013, MNRAS, 432, 1483

Carrasco Kind M., Brunner R. J., 2014, MNRAS, 441, 3550

Chang C., et al., 2013, MNRAS, 434, 2121

Choromanska A., Henaff M., Mathieu M., Ben Arous G., LeCun Y., 2014, arXiv e-prints, p. arXiv:1412.0233

Choudhury S. R., Choubey S., 2018, Journal of Cosmology and Astro-Particle Physics, 2018, 017

Coe D., Benítez N., Sánchez S. F., Jee M., Bouwens R., Ford H., 2006, AJ, 132, 926

Collister A. A., Lahav O., 2004, PASP, 116, 345

DES Collaboration et al., 2018a, arXiv e-prints, p

DES Collaboration et al., 2018b, arXiv e-prints, p. arXiv:1811.02375

D'Isanto A., Polsterer K. L., 2018, A&A, 609, A111

Dark Energy Survey Collaboration et al., 2016, MNRAS, 460,

Dawson W., Schneider M., 2014, Technical report, Complementarity of LSST and WFIRST: Regarding Object Blending. Lawrence Livermore National Laboratory (LLNL), Livermore, CA

Dempster A. P., Laird N. M., Rubin D. B., 1977, Journal of the Royal Statistical Society: Series B (Methodological), 39, 1

Duncan K. J., Jarvis M. J., Brown M. J., Röttgering H. J., 2018, Monthly Notices of the Royal Astronomical Society, 477, 5177

Edge A., Sutherland W., Kuijken K., Driver S., McMahon R., Eales S., Emerson J. P., 2013, The Messenger, 154, 32

Feroz F., Hobson M. P., Bridges M., 2009, MNRAS, 398, 1601 Foreman-Mackey D., 2016, The Journal of Open Source Software, 24

Foreman-Mackey D., Hogg D. W., Lang D., Goodman J., 2013, PASP, 125, 306

Gelman A., Carlin J. B., Stern H. S., Dunson D. B., Vehtari A., Rubin D. B., 2013, Bayesian data analysis. Chapman and Hall/CRC

Gerdes D. W., Sypniewski A. J., McKay T. A., Hao J., Weis M. R., Wechsler R. H., Busha M. T., 2010, ApJ, 715, 823

Gomes Z., Jarvis M. J., Almosallam I. A., Roberts S. J., 2017, Monthly Notices of the Royal Astronomical Society, 475, 331

Graham M. L., Connolly A. J., Ivezić Ž., Schmidt S. J., Jones R. L., Jurić M., Daniel S. F., Yoachim P., 2018, AJ, 155, 1

Hildebrandt H., et al., 2010, A&A, 523, A31

Hill D. T., et al., 2011, MNRAS, 412, 765

 $Hoekstra\ H.,\ et\ al.,\ 2006,\ The\ Astrophysical\ Journal,\ 647,\ 116$

```
Hoerl A. E., Kennard R. W., 1970, Technometrics, 12, 55
```

Holwerda B. W., et al., 2015, MNRAS, 449, 4277

Ilbert O., et al., 2006, A&A, 457, 841

Ivezić Ž., et al., 2019, ApJ, 873, 111

Jones D. M., Heavens A. F., 2019, MNRAS, 483, 2487

Joseph R., Courbin F., Starck J.-L., 2016, A&A, 589, A2

Kitzbichler M. G., White S. D. M., 2007, MNRAS, 376, 2

Kremer J., Gieseke F., Steenstrup Pedersen K., Igel C., 2015, Astronomy and Computing, 12, 67

LSST Science Collaboration et al., 2009, preprint, (arXiv:0912.0201)

Laureijs R., et al., 2011, preprint, (arXiv:1110.3193)

Leistedt B., Hogg D. W., 2017, The Astrophysical Journal, 838,

Lima M., Cunha C. E., Oyaizu H., Frieman J., Lin H., Sheldon E. S., 2008, MNRAS, 390, 118

Lupton R. H., 2005, Technical report, SDSS image processing I: The deblender

MacKay D. J., 1992, Neural computation, 4, 448

Malz A. I., Marshall P. J., DeRose J., Graham M. L., Schmidt S. J., and R. W., 2018, The Astronomical Journal, 156, 35

Mandelbaum R., et al., 2008, MNRAS, 386, 781 Melchior P., Moolekamp F., Jerdee M., Armstrong R., Sun A.-L.,

Bosch J., Lupton R., 2018, preprint, (arXiv:1802.10157) Meshcheryakov A., Glazkova V., Gerasimov S., Mashechkin I., 2018, Astronomy Letters, 44, 735

Myers A. D., White M., Ball N. M., 2009, MNRAS, 399, 2279

Peebles P. J. E., 2001, in Martínez V. J., Trimble V., Pons-Bordería M. J., eds, Astronomical Society of the Pacific Conference Series Vol. 252, Historical Development of Modern Cosmology. p. 201 (arXiv:astro-ph/0103040)

Petersen K. B., Pedersen M. S., 2014

Prechelt L., 1998, Early stopping-but when?. Springer, doi:10.1007/3-540-49430-8_3

Racca G. D., et al., 2016, in Space Telescopes and Instrumentation 2016: Optical, Infrared, and Millimeter Wave. p. 99040O (arXiv:1610.05508), doi:10.1117/12.2230762

Rhodes J., et al., 2017, ApJS, 233, 21

Rivera J. D., Moraes B., Merson A. I., Jouvel S., Abdalla F. B., Abdalla M. C. B., 2018, MNRAS, 477, 4330

Rodríguez-Muñoz L., et al., 2019, Monthly Notices of the Royal Astronomical Society, 485, 586

Runnalls A. R., 2007, IEEE Transactions on Aerospace and Electronic Systems, 43, 989

Sadeh I., Abdalla F. B., Lahav O., 2016, PASP, 128, 104502 Schieferdecker D., Huber M. F., 2009, in 2009 12th International

Conference on Information Fusion. pp 1536–1543
Sellentin E., Starck J.-L., 2019, arXiv e-prints, p.

arXiv:1902.00709
Shrikumar A., Greenside P., Kundaje A., 2017, arXiv e-prints, p.

arXiv:1704.02685

Shwartz-Ziv R., Tishby N., 2017, arXiv e-prints, p. arXiv:1703.00810

Skilling J., 2006, Bayesian Anal., 1, 833

Soo J. Y. H., et al., 2018, MNRAS, 475, 3613

Stoughton C., et al., 2002, AJ, 123, 485

Syring N., Martin R., 2018, Biometrika, 106, 479

Wadadekar Y., 2005, PASP, 117, 79

Way M. J., Srivastava A. N., 2006, ApJ, 647, 102

West M., 1993, Journal of the Royal Statistical Society: Series B (Methodological), 55, 409

Williams J. L., Maybeck P. S., 2006, Mathematical and Computer Modelling, 43, 976

Wittman D., Bhaskar R., Tobin R., 2016, MNRAS, 457, 4005

This paper has been typeset from a TEX/LATEX file prepared by the author.