



Statistical Modeling of Galaxy Catalogues with Normalizing Flows

2 JOHN FRANKLIN CRENSHAW,^{1,2} J. BRYCE KALMBACH,¹ ALEXANDER GAGLIANO,^{3,4,5,6} ZIANG YAN,⁷ ANDREW J. CONNOLLY,¹
3
4 THE LSST DARK ENERGY SCIENCE COLLABORATION

6 ¹*DIRAC Institute and the Department of Astronomy, University of Washington, Seattle, WA 98195, USA*

7 ²*Department of Physics, University of Washington, Seattle, WA 98195, USA*

8 ³*Department of Astronomy, University of Illinois at Urbana-Champaign, 1002 W. Green St., IL 61801, USA*

9 ⁴*National Center for Supercomputing Applications, Urbana, IL, 61801, USA*

10 ⁵*Center for AstroPhysical Surveys, Urbana, IL, 61801, USA*

11 ⁶*National Science Foundation Graduate Research Fellow*

12 ⁷*Ruhr University Bochum, Faculty of Physics and Astronomy, Astronomical Institute (AIRUB), German Centre for Cosmological
13 Lensing, 44780 Bochum, Germany*

14 ⁸*McWilliams Center for Cosmology, Department of Physics, Carnegie Mellon University*

15 ⁹*Department of Physics and Astronomy, University of California, One Shields Avenue, Davis, CA 95616, USA*

ABSTRACT

17 Normalizing flows are powerful tools for learning high-dimensional probability distributions from
18 samples thereof, and have many applications in astronomy, including posterior estimation and forward
19 modeling. We introduce PZFlow, a Python package for statistical modeling of tabular data using nor-
20 malizing flows. We use PZFlow to model photometric galaxy catalogs including redshifts, photometry,
21 size, and shape, and generate a synthetic catalog. In this catalog, each galaxy has a *true* redshift pos-
22 terior. We discuss the importance of these true posteriors for the comprehensive evaluation of redshift
23 posteriors produced by photometric redshift (photo-z) estimators. We also demonstrate the use of an
24 ensemble of normalizing flows for density estimation, applied to photo-z estimation. While we focus
25 on photo-z estimation and validation, we emphasize that these methods are applicable to any galaxy
26 properties, and any other tabular data.

1. INTRODUCTION

27 Astronomical data represent realizations of complex
28 probability distributions. A common goal in research is
29 to infer the underlying distribution from a limited set
30 of noisy data, and use this distribution to estimate pos-
31 terior distributions over galaxy and population param-
32 eters. An important example is photometric redshift
33 (photo-z) estimation, where galaxy redshift posteriors
34 are estimated from galaxy photometry, using a model
35 informed by a training set of spectroscopic galaxy data
36 ([Newman & Gruen 2022](#)). These redshift posteriors are
37 then used to estimate posterior distributions for cosmo-
38 logical parameters ([The LSST Dark Energy Science Col-](#)
39 [laboration et al. 2018](#)).

41 Knowledge of the probability distributions underly-
42 ing our data is also valuable for simulation and forward
43 modeling astronomical data sets. Simulating data sets
44 with realistic statistical properties enables methodolog-
45 ical development and the calibration of systematic un-
46 certainties. Forward modeling also facilitates data aug-
47 mentation (e.g. [Lokken et al. 2022](#)) and the creation
48 of large data challenges (e.g. [Kessler et al. 2019; LSST](#)
49 [Dark Energy Science Collaboration et al. 2021; Korytov](#)
50 [et al. 2019](#)), which are becoming more prevalent in the
51 big data era of astronomy.

52 Inferring the underlying probability distribution, or
53 likelihood, from a high-dimensional data set is a diffi-
54 cult problem. Many machine learning tools have been
55 developed for this task, including Generative Adversar-
56 ial Networks (GANs; [Goodfellow et al. 2014](#)) and Vari-
57 ational Autoencoders (VAEs; [Kingma & Welling 2014](#)).
58 Both are neural networks that excel at forward mod-
59 eling complex data sets, but neither allow exact likeli-

Corresponding author: John Franklin Crenshaw
jfc20@uw.edu

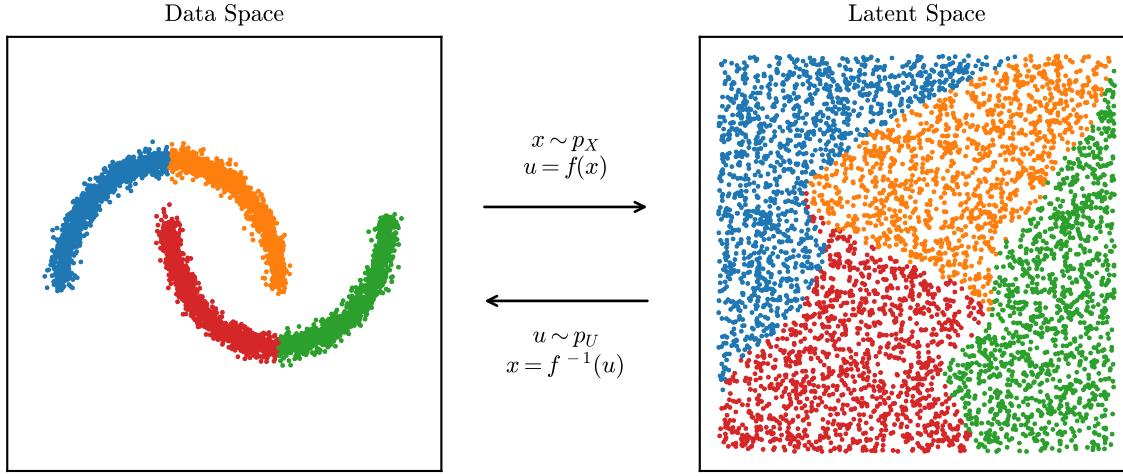


Figure 1. A normalizing flow demonstrated on the two moons data set from scikit-learn. The two moons data on the left is mapped onto a two dimensional uniform distribution by the bijection f . The data are colored by quadrant to visualize their image in the latent space. You can sample the data distribution by sampling from the uniform distribution, and using f^{-1} to map the samples back to the data space.

hood calculation for the simulated data. This means that questions such as “under my model, what is the posterior distribution for redshift given the simulated photometry” have only approximate answers.

Normalizing flows, on the other hand, are a deep learning tool that excel at forward modeling, while also allowing exact, analytic likelihood calculation. In other words, they allow you to analytically calculate likelihoods and posteriors with respect to the distribution from which simulated data is drawn. Normalizing flows operate by learning an invertible transformation of the data distribution into a simpler, tractable distribution, known as the latent distribution. A common choice for latent distribution is a normal distribution, hence the name *normalizing* flow. This allows sampling and likelihood calculation to be performed within the context of the simple latent distribution, with the normalizing flow acting as a translator between the latent samples and likelihoods, and their corresponding values in the context of the complicated data distribution.

Normalizing flows have gained popularity as tools for efficient and flexible sampling for parameter inference (e.g., Dai & Seljak 2022; Dacunha et al. 2022; Hassan et al. 2022). In this paper, we focus on photo-z’s and forward modeling photometric galaxy catalogs. Since normalizing flows allow exact, analytic probability calculation, the properties of objects in these catalogs have *true* posteriors. This means that questions such as “under my model, what is the posterior distribution for redshift given the simulated photometry” have exact answers. Catalogs with true posteriors are useful for the testing and validation of analysis pipelines that estimate posteriors, such as the photo-z estimators used in much of as-

trophysics and cosmology. Previous evaluations of these estimators have focused on comparing point estimates to true values (e.g., Hildebrandt et al. 2010; Sánchez et al. 2014; Graham et al. 2018), or evaluating ensembles of posteriors (Schmidt & Malz et al. 2020). Normalizing flow catalogs with true posteriors open up a new, more comprehensive avenue for evaluation of these estimators by enabling direct posterior-to-posterior comparison.

To facilitate the statistical modeling of galaxy catalogs and other astronomical data sets, we have developed PZFlow, a normalizing flow package for Python. With relatively little tuning required by the user, PZFlow can provide a generative model for any tabular data, including continuous and discrete variables, and variables with Euclidean or periodic topology. In addition to generative modeling, PZFlow can calculate posteriors over any columns in your data set, and can convolve errors and marginalize over missing columns while training the model or calculating posteriors for samples.

In this paper, we provide the background on normalizing flows (Section 2) required to understand PZFlow (Section 3). We then use PZFlow to simulate a galaxy catalog (Section 4), where each object has photometry, size, ellipticity, redshift, and a true redshift posterior. We also demonstrate using PZFlow as a density estimator, via the example of photo-z estimation (Section 5). We conclude in Section 6.

2. NORMALIZING FLOWS

Normalizing flows model complex, high-dimensional probability distributions by learning a mapping from the

123 data distribution to a tractable latent distribution¹. Of
 124 ten the latent distribution is a standard Normal distribu-
 125 tion, and so the mapping “normalizes” the data, hence
 126 the name “normalizing flow”. This mapping allows us
 127 to sample and evaluate densities using the latent distri-
 128 bution, rather than the unknown data distribution.

129 Assume we have a differentiable function f that maps
 130 samples x from the data distribution p_X onto samples
 131 u from the latent distribution p_U . Using the change of
 132 variables formula, we can evaluate the probability den-
 133 sity of the data:

$$134 \quad p_X(x) = p_U(u = f(x)) |\det \nabla f(x)|, \quad (1)$$

136 where $\nabla f(x)$ is the Jacobian of f evaluated at x . In
 137 words, computing the density $p_X(x)$ is accomplished by
 138 mapping x to the latent distribution, calculating its den-
 139 sity there, and multiplying by the associated Jacobian
 140 determinant, which accounts for how the function f dis-
 141 torts volume elements of the space.

142 If we further assume that f is invertible, we can sample
 143 from the data distribution by applying f^{-1} to samples
 144 from the latent distribution²:

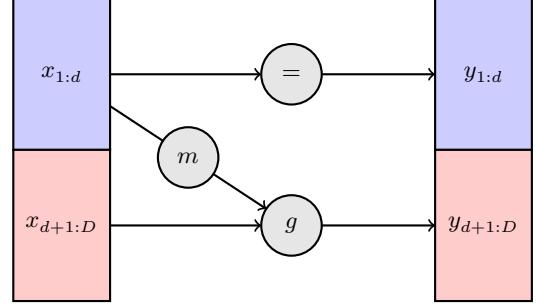
$$145 \quad x = f^{-1}(u) \quad \text{where} \quad u \sim p_U. \quad (2)$$

147 Figure 1 shows an example of a normalizing flow that
 148 transforms the scikit-learn (Pedregosa et al. 2011) two
 149 moons distribution into a uniform distribution. The
 150 data points are colored by quadrant to visualize their
 151 image under f .

152 The following sections discuss how to build a normal-
 153 izing flow to model data with various features. Section
 154 2.1 discusses the bijection f and introduces the building
 155 blocks from which our bijections will be built; Section
 156 2.2 discusses how to choose an appropriate latent distri-
 157 bution for your data; Section 2.3 describes how to build
 158 a flow that models a conditional distribution; Section
 159 2.4 explains how to model data with periodic topology;
 160 finally Section 2.5 explains how to model data with dis-
 161 crete variables.

162 2.1. Designing a bijection

163 A bijection is an invertible map between two sets. In
 164 a normalizing flow, the bijection maps the data distribu-
 165 tion onto the latent distribution for likelihood calcula-
 166 tion, and the inverse of the bijection maps samples from
 167 the latent distribution back to the data distribution.
 168 The bijection of a normalizing flow must be powerful



197 **Figure 2.** Diagram of a coupling layer. The first partition,
 198 $x_{1:d}$, is passed through the layer unchanged. The second par-
 199 tition, $x_{d+1:D}$, is transformed by the coupling law g , which
 200 is parameterized by the coupling function m applied to the
 201 first partition.

199 enough to model complex relationships in data, while
 200 remaining invertible and simultaneously possessing an
 201 efficiently computable Jacobian determinant. This lat-
 202 ter constraint is the primary difficulty in designing a nor-
 203 malizing flow. The most popular strategy for achieving
 204 these requirements is to exploit the fact that a composi-
 205 tion of bijections is also bijective. By chaining together
 206 multiple less-expressive bijections whose Jacobians are
 207 efficiently computable, a composite bijections can be
 208 constructed that meets our requirements:

$$179 \quad f = \cdots \circ f_3 \circ f_2 \circ f_1. \quad (3)$$

181 The overall Jacobian determinant can be efficiently cal-
 182 culated using the chain rule.

183 There is an extensive literature on constructing these
 184 sub-bijections (see Kobyzev et al. 2020 for a review).
 185 Some bijections are specialized to be particularly effi-
 186 cient at either density estimation or sampling, but for
 187 many science cases, we wish to do both. For this rea-
 188 son, we will focus on Rational-Quadratic Rolling Spline
 189 Couplings (RQ-RSCs), bijections which achieve state-
 190 of-the-art performance, while being efficient with both
 191 tasks (Durkan et al. 2019).

192 2.1.1. Rational-Quadratic Rolling Spline Couplings

193 RQ-RSCs are bijections that are composed of coupling
 194 layers (Dinh et al. 2015, 2017). A coupling layer parti-
 195 tions the data, $x \in \mathbb{R}^D$, into two sets, $x_{1:d}$ and $x_{d+1:D}$.
 196 The first set is then used to transform the second set:

$$197 \quad y_{1:d} = x_{1:d} \\ 198 \quad y_{d+1:D} = g(x_{d+1:D}; m(x_{1:d})), \quad (4)$$

199 where $g : \mathbb{R}^{D-d} \times \mathbb{R}^d \rightarrow \mathbb{R}^{D-d}$ is an invertible *coupling*
 200 *law*, and m is a *coupling function* defined on \mathbb{R}^d . This is
 201 illustrated in Figure 2. The advantage of this structure

¹ Some of the machine learning literature defines the mapping in the opposite direction.

² Here, \sim means “is drawn from.”

is that the Jacobian is triangular,

$$\frac{\partial y}{\partial x} = \begin{pmatrix} I_d & 0 \\ \frac{\partial y_{d+1:D}}{\partial x_{1:d}} & \frac{\partial y_{d+1:D}}{\partial x_{d+1:D}} \end{pmatrix}, \quad (5)$$

where I_d is the $d \times d$ identity matrix. In particular, the Jacobian determinant is

$$\det \frac{\partial y}{\partial x} = \det \frac{\partial y_{d+1:D}}{\partial x_{d+1:D}}. \quad (6)$$

Furthermore, the inverse can be calculated as

$$\begin{aligned} x_{1:d} &= y_{1:d} \\ x_{d+1:D} &= g^{-1}(y_{d+1:D}; m(x_{1:d})), \end{aligned} \quad (7)$$

Notice that neither inverting a coupling layer g , nor calculating the Jacobian determinant, requires inverting or taking derivatives of the coupling function m , which can thus be arbitrarily complex.

The obvious limitation of a coupling layer is that only a subset of the data dimensions are transformed. This is overcome by stacking multiple coupling layers in succession, and switching which variables belong to which partition. In practice, this is achieved by interspersing coupling layers with bijections that shuffle the dimensions of x . These shuffling bijections are trivially inverted and have a Jacobian determinant of one.

In a general coupling layer g , there are a variety of coupling laws m one can use. RQ-RSC's use Rational-Quadratic Neural Spline Coupling (Durkan et al. 2019). As the name suggests, the coupling law g is a set of rational-quadratic splines. In particular, $g_i : [-B, B] \rightarrow [-B, B]$ for each dimension i of $x_{d+1:D}$, where g_i is a piecewise combination of K segments, and each segment is a rational-quadratic function. The positions and derivatives of the knots that parameterize the splines are calculated using the coupling function m , which is a dense neural network applied to $x_{1:d}$.

The result is a bijection that achieves state-of-the-art performance and efficiency for forward modeling and density estimation (Kobyzev et al. 2020), and are flexible enough to model complex distributions with multiple discontinuities and hundreds of modes. In addition, they are easily adaptable for flows with periodic topology (Section 2.4). For more details, see Durkan et al. (2019).

In this work, we stack Rational-Quadratic Neural Spline Couplings, with Rolling Layers between each – a configuration we name Rational-Quadratic Rolling Spline Couplings (RQ-RSCs). Rolling Layers shift the dimensions of x by one place:

$$\text{Roll} : [x_1, \dots, x_{D-1}, x_D] \rightarrow [x_D, x_1, \dots, x_{D-1}]. \quad (8)$$

By constructing a stack with D coupling layers, RQ-RSCs individually transform each of the D dimensions of x as a function of the other $D - 1$ dimensions. In the limit of high spline resolution (i.e. $K \rightarrow \infty$), RQ-RSCs can model any differentiable, monotonic function on $[-B, B]^D$ and can thus model arbitrarily complex distributions in this region. In practice, we find very good performance for diverse data sets with $K \approx 16$.

Note you can specify a different value of K for each of the D spline layers in order to individually control the resolution of each dimension. Lowering K typically results in a smoother distribution, while increasing K increases the complexity the normalizing flow can capture, while also increasing computational and memory cost.

2.1.2. Data processing bijections

While RQ-RSCs perform the heavy lifting of mapping the data distribution p_X onto the latent distribution p_U , it is also convenient to define other bijections that perform useful operations such as pre- and post-processing. We name these *data processing bijections*.

For example, RQ-RSCs (and the RQ-NSCs on which they are based) are defined on the domain $[-B, B]$, and thus will not transform samples outside this range. It is therefore useful to define a *Shift Bounds* bijection, which shifts the original range of each dimension to match the domain of the splines. Note this shift must be set at training time, with the assumption that future test data will lie within the same bounds³. You can choose a range wider than that covered by the training set if you wish to allow the flow to sample outside the range of the training set

For an example of building an application-specific data processing bijection, see the *Color Transform* bijection defined in Section 4.1, which maps galaxy magnitudes to galaxy colors. See section 2.5 for data processing bijections that enable modeling of discrete data.

Instead of using these data processing bijections, you can of course manually pre-process the data before evaluating densities and post-process samples drawn from the normalizing flow. However, by building pre- and post-processing directly into the bijection, you remove these extra steps from the workflow. This reduces the complexity of working with the normalizing flow and ensures that the flow always “remembers” how to correctly perform these pre- and post-processing steps.

³ While this sounds quite restrictive, neural networks are typically pretty bad at extrapolating beyond the bounds of the training set anyway.

296 **2.2. Choosing a latent distribution**

297 In principle, with a sufficiently expressive bijection,
 298 the choice of latent distribution does not matter as long
 299 as it is a distribution in which you can easily sample
 300 and calculate densities. However, in practice, bijections
 301 are limited in expressiveness, i.e. they cannot necessarily
 302 transform any arbitrary data distribution into any
 303 arbitrary latent distribution.

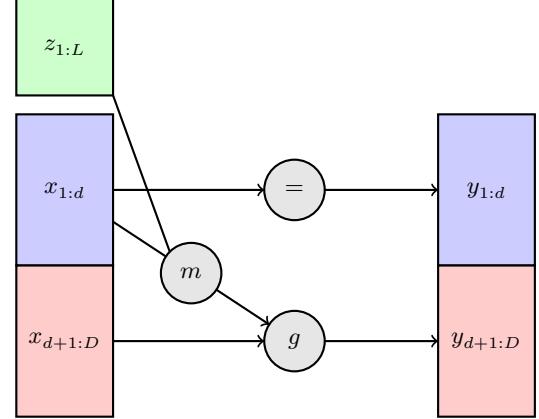
304 For example, the splines of RQ-RSCs only transform
 305 samples in the range $[-B, B]$. Sampling from a latent dis-
 306 tribution with support outside this range will therefore
 307 result in strange outliers and incorrect boundary con-
 308 ditions. One can apply a transformation to the latent
 309 samples before they are fed into the RQ-RSC to ensure
 310 that they lie within the support of the splines, but it is
 311 simpler to use a compact latent distribution whose sup-
 312 port matches that of the spline layers. A simple choice
 313 would be the uniform distribution over $[-B, B]$.

314 Additionally, as no bijection is perfect, the structure of
 315 the latent distribution will not be completely erased in
 316 the translation from latent to data distribution. Thus,
 317 the latent distribution can be viewed as a prior or induc-
 318 tive bias on samples from the data distribution (Jaini
 319 et al. 2020). It is therefore advantageous to select a
 320 latent distribution whose features match some of the
 321 structure in the data.

322 A latent distribution that can achieve both desider-
 323 ata is the Beta distribution, i.e. $u \sim \text{Beta}(\alpha, \beta)$, where
 324 $\alpha, \beta > 0$ are learnable parameters⁴. This distribution is
 325 compact, and by varying α and β this distribution can
 326 take on a wide variety of shapes with different means,
 327 skews, and kurtoses, allowing the inductive bias of the
 328 prior to adapt to structure in the data during train-
 329 ing. However, as RQ-RSCs are defined on the domain
 330 $[-B, B]$, it is more convenient to use a modified Beta
 331 distribution, which we name the *Centered Beta distri-*
 332 *bution*:

$$333 \quad \text{CentBeta}(u|\alpha, \beta, B) = 2B \left(\text{Beta}(u|\alpha, \beta) - \frac{1}{2} \right). \quad (9)$$

335 In general, as long as sampling and density evalua-
 336 tion are tractable, one can use any parameterization of
 337 the latent distribution that matches some desired struc-
 338 ture in the data and learn the distribution parameters
 339 during training. We give this generalization the name
 340 *latent-adaptive flows* (LAFs; inspired by the Tail Adap-
 341 tive Flows of Jaini et al. 2020). Our experiments in-
 342 dicate that learnable latent distributions can improve
 343 training loss, but require more care in training.



344 **Figure 3.** Diagram of a *conditional* coupling layer. The
 345 first partition, $x_{1:d}$, is passed through the layer unchanged.
 346 The second partition, $x_{d+1:D}$, is transformed by the coupling
 347 law g , which is parameterized by the coupling function m
 348 applied to the first partition *and* the conditional variables
 349 $z_{1:L}$. The conditional variables are *never* altered by the flow.

354 Note that while we discussed univariate distributions
 355 above, these considerations generalize easily to multi-
 356 ple dimensions. Each of these distributions have multi-
 357 variate generalizations that can be used when modeling
 358 higher-dimensional data. The full multivariate latent
 359 distribution can also be assembled by taking the prod-
 360 uct of multiple univariate distributions⁵. This may even
 361 be desired if different dimensions of the data have dif-
 362 ferent structure that you wish to encode in the latent
 363 distribution.

354 **2.3. Conditional flows**

355 The bijections and latent distributions discussed
 356 above can be easily adapted to directly learn conditional
 357 probability distributions: you only need to make the re-
 358 placement $f(x) \rightarrow f(x; z)$, where z is a vector of condi-
 359 tions (Winkler et al. 2019). This is illustrated in Figure
 360 3, which is a modification of Figure 2 to include the
 361 input of conditional variables to the coupling function
 362 m . In practice, since m is usually a neural network,
 363 this amounts to just appending the conditions z to the
 364 inputs of the neural network.

365 While $p(x|z)$ is technically encoded within $p(x, z)$,
 366 which can be learned with a regular normalizing flow,
 367 directly modeling $p(x|z)$ with a conditional flow has a
 368 few benefits. Training is typically faster, since the latent
 369 distribution has a smaller number of dimensions. You
 370 can also draw samples of x at fixed values of the condi-
 371 tions z , and you can calculate $p(x|z)$ without having to

⁴ In practice, it is easier to learn $\log \alpha$ and $\log \beta$ to ensure that $\alpha, \beta > 0$.

⁵ Note that while the latent variables will be independent, the data variables will still have correlations imprinted by the bijections.

372 numerically calculate and divide by $p(z)$, which can be
 373 computationally expensive.

374 2.4. Flows with periodic topology

375 The flows we have considered so far model data that
 376 live in \mathbb{R}^n . This assumption is insufficient for modeling
 377 variables from spaces with non-Euclidean topology, e.g.
 378 positions on the sky. While progress has been made on
 379 building flows for general topologies (e.g. [Gemici et al. 2016](#) and [Falorsi et al. 2019](#)), we will focus on building
 380 flows on the sphere, S^2 , as this is the case most relevant
 381 in astronomy. We will see that by carefully choosing
 382 the latent space, we can construct flows with periodic
 383 topology with minimal additional effort ([Rezende et al. 2020](#)).
 384

385 Positions on the sphere are specified by two angles⁶,
 386 θ and ϕ , the latter of which is periodic. By mapping
 387 θ to $\cos\theta$, we map the sphere to a cylinder⁷: $S^2 \rightarrow$
 388 $[-1, 1] \times S^1$ (i.e. the Cartesian product of an interval
 389 and a circle). In other words, we can transform $\cos\theta$
 390 with a Euclidean flow, as long as we ensure that the flow
 391 bounds samples to the range $[-1, 1]$. However, the S^1
 392 piece, ϕ , has a periodic topology and must be handled
 393 more carefully.

394 First, we will address transformations of $\cos\theta$. The
 395 only constraint we must impose is that samples of $\cos\theta$
 396 must lie in the range $[-1, 1]$. Fortunately, RQ-RSCs are
 397 bounded, mapping a range in u to the same range in
 398 x . Thus, if we pick a latent distribution with compact
 399 support in $[-1, 1]$, samples of $\cos\theta$ are guaranteed to lie
 400 in the same range, as long as we set the range of the
 401 RQ-RSC $B = 1$.

402 Next we will address transformations of ϕ . For f to be
 403 a differentiable bijection on the circle, S^1 , it is sufficient
 404 that f obey the following constraints:

$$406 \quad f(0) = 0 \quad (10)$$

$$407 \quad f(2\pi) = 2\pi \quad (11)$$

$$408 \quad \nabla f(0) = \nabla f(2\pi) \quad (12)$$

$$409 \quad \nabla f(\phi) > 0. \quad (13)$$

410 The first two constraints ensure continuity of f by des-
 411 ignating $\phi = 0$ as a fixed point, and the third constraint
 412 ensures continuity of ∇f at that fixed point. While the
 413 designation of $\phi = 0$ as a fixed point is an unnecessary
 414 restriction on f , any differentiable bijection on the circle
 415

⁶ We use the convention where θ and ϕ are the zenith and azimuthal angles, respectively.

⁷ This map can be explicitly constructed via an embedding in \mathbb{R}^3 . Technically, the map is not defined for $\theta \in \{0, \pi\}$, however as this set has zero measure, it can be safely ignored.

416 has at least one fixed point up to a phase change, and so
 417 this restriction does not actually restrict the expressive-
 418 ness of f . The fourth restriction ensures monotonicity,
 419 which guarantees invertibility.

420 If we make the phase change $\phi \rightarrow \phi - \pi$ so that our
 421 angles $\phi \in [-\pi, \pi]$, a RQ-NSC with $B = \pi$ automatically
 422 fulfills all four constraints. In fact, regular RQ-NSC's
 423 impose the further condition

$$424 \quad \nabla f(-\pi) = \nabla f(\pi) = 1 \quad (14)$$

425 to match an identity transform for inputs outside of the
 426 range $[-\pi, \pi]$. By choosing a latent distribution with
 427 compact support in the range $[-\pi, \pi]$, we ensure that
 428 no samples will lie outside the range of the splines, and
 429 so we can relax the boundary condition of Equation
 430 14 in favor of the boundary condition in Equation 12.
 431 Spline transforms with this relaxed boundary condition
 432 are named *Circular Splines* by [Rezende et al. \(2020\)](#).

433 The circular spline construction above is easily gener-
 434 alized to n-spheres and n-tori: $S^n \rightarrow [-1, 1]^{n-1} \times S^1$
 435 and $T^n \rightarrow (S^1)^n$ (see [Rezende et al. 2020](#) for more
 436 details). We can model the joint distribution of peri-
 437 odic and non-periodic variables with RQ-RSCs simply
 438 by choosing appropriate bounds B for each dimension,
 439 and by swapping boundary condition 14 for condition
 440 12 for any periodic dimensions.

442 2.5. Modeling discrete variables

443 In addition to the continuous variables described
 444 above, normalizing flows can also be used to model dis-
 445 crete variables. This can be achieved by “dequantizing”
 446 the discrete dimensions of the data, which can then be
 447 mapped onto continuous latent distributions using reg-
 448 ular continuous bijections. Dequantization consists of
 449 adding some kind of continuous noise to the discrete
 450 dimensions, transforming them into continuous dimen-
 451 sions. When sampling from the flow, you simply do the
 452 opposite, and “quantize” the discrete dimensions after
 453 applying all of the regular bijections, mapping the noisy,
 454 continuous variables onto their discrete counterparts.

455 A common method for dequantization is uniform de-
 456 quantization, in which random uniform noise in the
 457 range $(0, 1)$ is added to the discrete dimensions. The
 458 corresponding quantization applied while sampling from
 459 the flow consists of applying the floor function to the de-
 460 quantized dimensions, mapping these samples onto the
 461 nearest integer less than the sampled value. More so-
 462 phisticated dequantization schemes use variational infer-
 463 ence or even another normalizing flow to determine the
 464 noise distributions, which improves results by smoothing
 465 the discontinuities between neighboring discrete values.
 466 See [Ho et al. \(2019\)](#) [Hoogeboom et al. \(2020\)](#) for more
 467 details.

468 While the dequantizers are not technically bijections,
 469 they can be treated as data processing bijections and
 470 be chained together with the other bijections in your
 471 normalizing flow.

472 3. PZFLOW

473 PZFlow is a Python package for building normalizing
 474 flows, with a focus on easy high-performance modeling
 475 of high-dimensional tabular data. Data is handled in
 476 Pandas DataFrames (Wes McKinney 2010), while the
 477 normalizing flows are implemented in Jax (Bradbury
 478 et al. 2018), which allows for efficient, parallelizable,
 479 GPU-enabled calculations for very large data sets. The
 480 code is easily installable from the Python Package In-
 481 dex⁸ (PyPI) and is hosted on Github.⁹. The documen-
 482 tation¹⁰ includes tutorial notebooks demonstrating the
 483 features mentioned in this paper on different example
 484 problems.

485 The rest of this paper will demonstrate using PZFlow
 486 for the statistical modeling of galaxy catalogs. Section
 487 4 uses PZFlow to forward model a galaxy catalog, in-
 488 cluding photometry, spectroscopic redshifts (spec-z’s),
 489 *true* photo-z posteriors, ellipticities, and sizes. Section
 490 5 uses PZFlow for photo-z estimation, demonstrating
 491 the power of PZFlow as a density estimator, including
 492 numerous useful features for photo-z estimation.

493 In addition to the examples in this paper, PZFlow has
 494 already been used in various other projects:

- 495 • Malz et al. (2021) used PZFlow to build a met-
 496 ric for observing strategy optimization based on
 497 information theory;
- 498 • Stylianou et al. (2022) used PZFlow to forward
 499 model galaxy data with true redshift posteriors in
 500 order to evaluate the impact of survey incomple-
 501 ness and spec-z errors on photo-z estimation;
- 502 • Lokken et al. (2022) used PZFlow to smooth high-
 503 redshift artifacts in simulations of host galaxies for
 504 supernovae and other transients.

505 4. FORWARD MODELING A GALAXY CATALOG

506 In this section, we use PZFlow to forward model a
 507 photometric galaxy catalog for the Vera Rubin Obser-
 508 vatory’s Legacy Survey of Space and Time (LSST; Ivezic
 509 et al. 2019). The advantage of using a catalog generated
 510 from a normalizing flow is that we have direct access to

511 the exact distribution from which the data is drawn, en-
 512 abling us to calculate true values for derived statistical
 513 products, such as the *true* photo-z redshift posterior for
 514 each galaxy.

515 In Section 4.1 we construct a normalizing flow to
 516 model the galaxy redshifts and photometry and gener-
 517 ate a new simulated catalog. In Section 4.2, we calculate
 518 true redshift posteriors for the new catalog. In Section
 519 4.3 we build a conditional flow to add additional galaxy
 520 properties to the catalog.

521 4.1. Forward modeling redshifts and photometry

522 To create a generative model of galaxy redshifts and
 523 photometry, we use the true redshifts and *ugrizy* magni-
 524 tudes from the CosmoDC2 simulation (LSST Dark En-
 525 ergy Science Collaboration et al. 2021; Korytov et al.
 526 2019) of the LSST Dark Energy Science Collaboration
 527 (DESC). We selected all galaxies from CosmoDC2 with
 528 at most one band with a signal-to-noise ratio (SNR) less
 529 than 10, using the forecasted photometric errors from
 530 LSST year 10. Of these, we randomly selected 10^6 galax-
 531 ies and split them into training and test sets consisting
 532 of 80% and 20% of the galaxies, respectively.

533 For the latent distribution we use a 7 dimensional Uni-
 534 form distribution over the range $[-5, 5]$ ¹¹. To map the
 535 data onto the latent distribution, we use the following
 536 bijection:

$$537 \quad f = \text{RQ-RSC} \circ \text{Shift Bounds} \circ \text{Color Transform}. \quad (15)$$

538 We will explain each layer of the bijection in the order
 539 they are applied to the input data.

The first layer of the bijection is the Color Transform, a data processing bijection designed specifically for this task. The Color Transform converts galaxy magnitudes to colors, but keeps the *i* band magnitude as a proxy for the apparent luminosity:

$$540 \quad \begin{aligned} \text{Color Transform : } & (\text{redshift}, u, g, r, i, z, y) \rightarrow \\ & (\text{redshift}, i, u - g, g - r, r - i, i - z, z - y). \end{aligned} \quad (16)$$

541 This layer is useful as galaxy redshifts correlate more
 542 directly with galaxy colors than galaxy magnitudes.

543 The next layer, Shift Bounds, is the data processing
 544 bijection defined in Section 2.1.2, which maps the range
 545 of the data onto the support of the RQ-RSC. Note that
 546 since Shift Bounds is on the “other side” of the Color
 547 Transform, we need to map the ranges of the colors $u - g$,
 548 $g - r$, etc. onto the support of the splines, instead of the
 549 original magnitudes.

⁸ <https://pypi.org/project/pzflow/>

⁹ <https://github.com/jfcrenshaw/pzflow>

¹⁰ <https://jfcrenshaw.github.io/pzflow/>

¹¹ The choice of 5 was arbitrary. Any other positive value would work just as well.

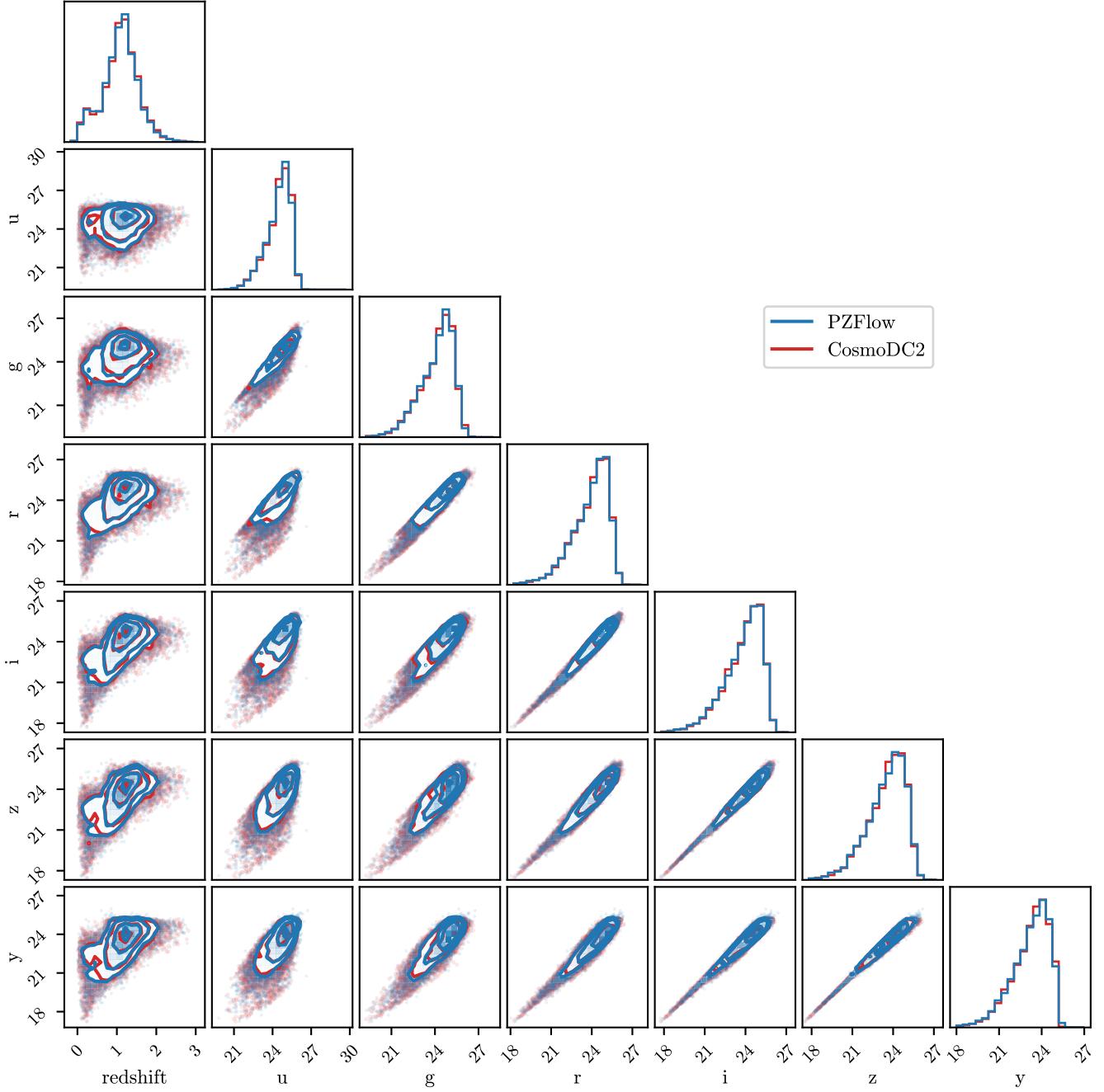


Figure 4. Distribution of true redshifts and noiseless photometry from the CosmoDC2 test set, compared to a sample drawn from the distribution learned by PZFlow. The close overlap of every pair-wise distribution demonstrates that PZFlow has learned the distribution in CosmoDC2 with high fidelity.

The final layer is an RQ-RSC, described in detail in Section 2.1.1. This layer performs the heavy lifting of transforming the data distribution into the uniform latent distribution. We use $D = 7$ layers to transform all 7 dimensions of our data, and set $B = 5$ to match the support of the latent distribution. We use the coupling function (a feedforward neural network with two hidden

layers of 128 neurons) described in Durkan et al. (2019).
We use $K = 16$ spline knots.

After training the flow (see Appendix A), we assess the results by drawing 10^4 galaxies from the trained flow, and plotting their distribution against 10^4 galaxies from the test set (Figure 4). We see the normalizing flow has done an excellent job of reproducing the distribution of galaxies in CosmoDC2, without any unusual artifacts or

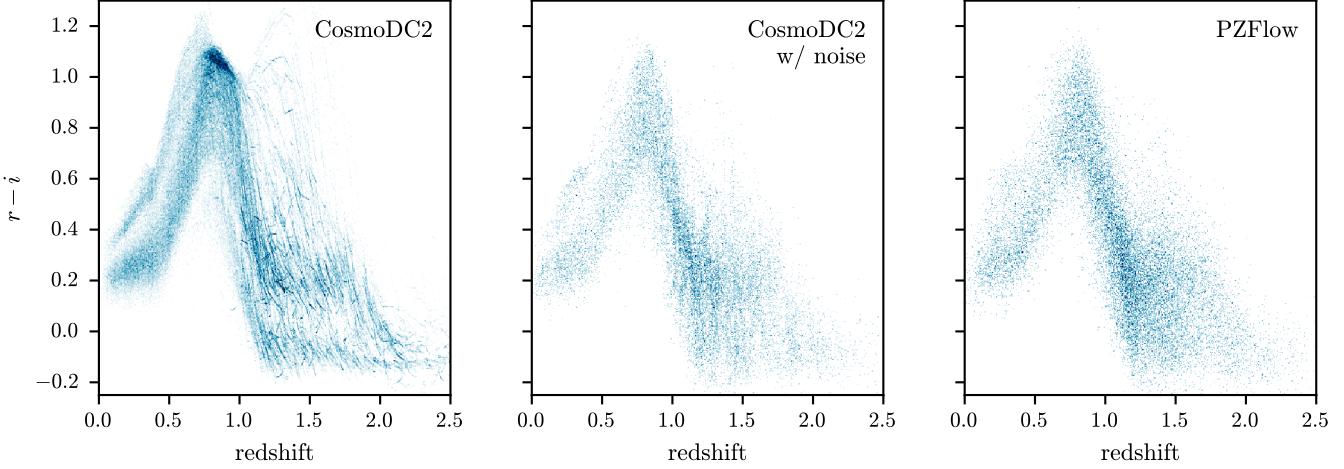


Figure 5. Comparing the noiseless $r - i$ vs redshift distribution for galaxy samples from CosmoDC2 (left) and from the normalizing flow (right). The high-redshift galaxies in CosmoDC2 lie along discrete tracks in color space due to the discrete number of galaxy SED templates used in the simulation. PZFlow smooths over these discrete tracks, resulting in a color distribution that is smooth up to high redshifts.

565 outliers. In addition, Figure 5 compares the distribution
 566 of galaxy $r - i$ vs redshift. The CosmoDC2 simulation
 567 is known to exhibit discrete tracks in this space at high
 568 redshift, due to the discrete number of SED templates
 569 used during simulation. These tracks are visible in the
 570 left panel. The right panel shows that PZFlow smooths
 571 over this discreteness, resulting in a color distribution
 572 that is smooth up to high redshifts.

573 We note that these results were obtained without any
 574 extensive hyperparameter search, and that very sim-
 575 ilar (slightly worse results) are obtained without the
 576 `ColorTransform` bijection, demonstrating the flexibility
 577 of the method to adapt to unseen data sets.

578 With this normalizing flow, we have an efficient Cos-
 579 moDC2 emulator that produces a smooth distribution
 580 of realistic galaxies up to high-redshifts. We use this
 581 emulator to generate a catalog with 10^4 galaxies. We
 582 add photometric errors using the LSST error model of
 583 our PhotErr package (see Appendix B). Importantly,
 584 since we have access to the probability distribution from
 585 which the galaxies were generated, we can calculate *true*
 586 redshift posteriors for each galaxy. This is the subject
 587 of the next section.

588 4.2. Calculating true posteriors

589 Since we have direct access to the probability distribu-
 590 tion from which the photometry and redshifts are drawn,
 591 we can analytically calculate the true redshift posterior
 592 for each galaxy: $p(z|m)$ where m is the vector of galaxy
 593 magnitudes. We note that this is not an estimate, like
 594 what would be returned by a photo-z estimator, but
 595 rather the truth, obtained from the model that gener-
 596 ated the photometry and redshifts in the first place.

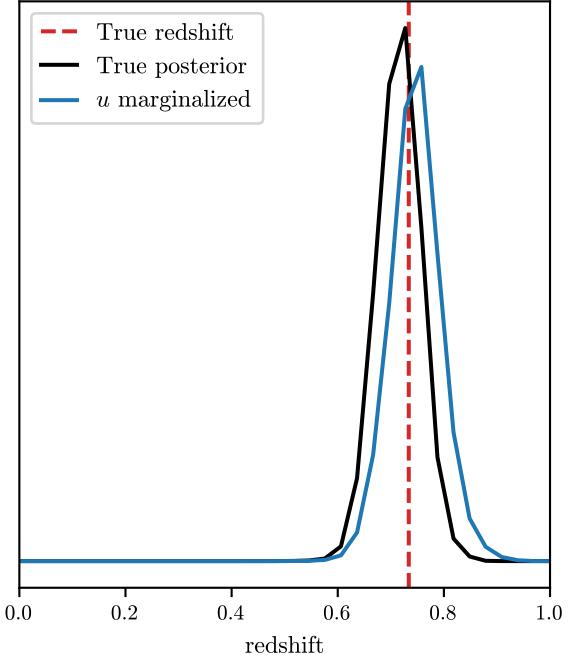


Figure 6. True redshift posteriors for a galaxy representing different amounts of information. The black posterior is calculated with the true magnitudes; the blue posterior is calculated after adding photometric errors; the orange posterior is calculated after adding photometric errors, but with the errors convolved during posterior calculation; the green posterior is the same as the orange, except with the u band marginalized over.

597 In addition to calculating the true posterior, $p(z|m)$,
 598 we can calculate a true posterior that is consistent with

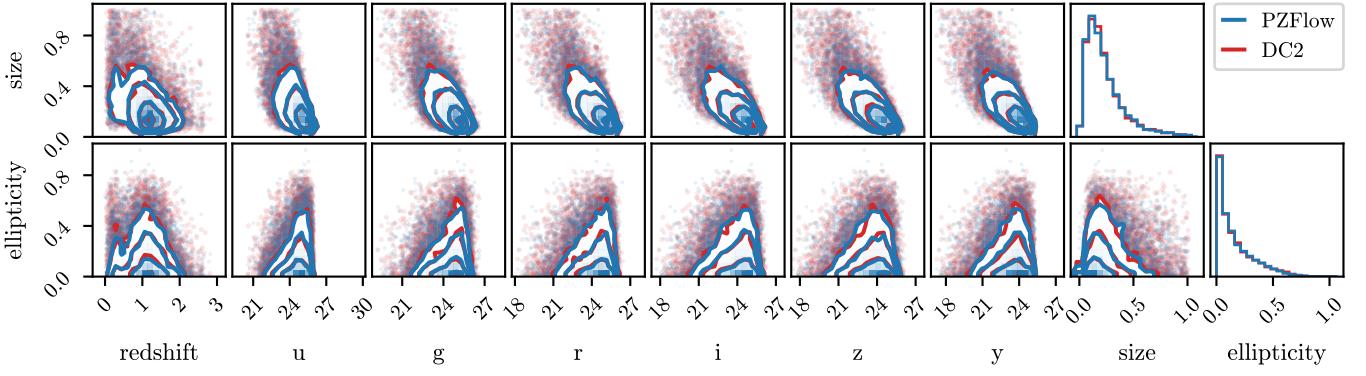


Figure 7. Conditional distributions of the ellipticity and size of the galaxies in the CosmoDC2 test set compared to the distribution learned by PZFlow. The close overlap of every pair-wise distribution demonstrates that PZFlow has learned the distribution in CosmoDC2 with high fidelity.

the photometric errors:

$$p(z|m, \sigma_m) = \int p(z|\hat{m})p(\hat{m}|m, \sigma_m)d\hat{m}, \quad (17)$$

where σ_m is the vector of photometric errors returned by the error model, and \hat{m} are possible noisy observations of the true magnitudes, m . We can evaluate this integral numerically by sampling $\hat{m} \sim p(\hat{m}|m, \sigma_m)$, evaluating $p(z, \hat{m})$ on a redshift grid, and normalizing to obtain $p(z|\hat{m})$. Averaging $p(z|\hat{m})$ over the samples \hat{m} yields $p(z|m, \sigma_m)$. This is the “best case scenario” redshift posterior that photo-z estimators hope to estimate.

We can also marginalize over any missing bands, n :

$$p(z|\hat{m}) = \frac{1}{p(\hat{m})} \int p(z, n, \hat{m})dn, \quad (18)$$

which can be calculated by evaluating $p(z, n, \hat{m})$ on a grid of z and n , summing over n to yield $p(z, \hat{m})$, and normalizing with respect to redshift to yield $p(z|\hat{m})$. You can once again average over \hat{m} samples to convolve the photometric errors. You may wish to marginalize over all values of n if the galaxy was not observed in that band. This might occur, for example, if you include photometry from the Euclid Space Telescope (Scaramella et al. 2022), which will not have complete coverage of the LSST catalog. You may also wish to marginalize over all values beyond a limiting magnitude if the galaxy was observed, but not detected in that band. This might occur, for example, in the low wavelength bands of Lyman dropout galaxies observed by LSST.

To visualize the impacts of error convolution and band marginalization, Figure 6 shows a number of redshift posteriors for a single example galaxy. The black posterior is calculated using the true galaxy magnitudes, while the blue posterior is calculated after adding photometric errors. Calculating the posterior using the

noisy photometry results in a biased posterior. The orange posterior has had the photometric errors convolved as in Equation 17. Convolving the errors broadens the posterior so that the true redshift lies within the support of the posterior. This broadening reflects the increased uncertainty due to the photometric errors. Finally, the green posterior has had the u band marginalized. This posterior favors higher redshifts, but still comfortably agrees with the true redshift. This demonstrates how the u band helps constrain the redshift, and the loss of this information leads to greater uncertainty.

Calculating these posteriors enables direct comparison of true redshift posteriors (consistent with photometric errors and missing bands) with the redshift posteriors estimated by photo-z estimators. This is important, as modern cosmology analyses are beginning to increasingly rely on full redshift posteriors (Mandelbaum et al. 2008; Newman & Gruen 2022). Schmidt & Malz (2020) showed that popular metrics for evaluating photo-z estimators using ensembles of photo-z posteriors can be misleading, and are not well suited to our science. PZFlow catalogs with *true* redshift posteriors provide a path forward by enabling the evaluation of photo-z estimators on a per-posterior basis.

4.3. Additional properties with conditional flows

In addition to the galaxy magnitudes and redshifts modeled above, we wish to include other galaxy properties in the catalog, such as galaxy size and ellipticity. In principle, we could have included these variables in the original normalizing flow. However, we did not want the true redshift posteriors to be conditioned on these variables, as most photo-z estimators only use galaxy photometry. Therefore, we will build a second flow that models these additional values conditioned on the galaxy redshift and magnitudes. Note that while we have only

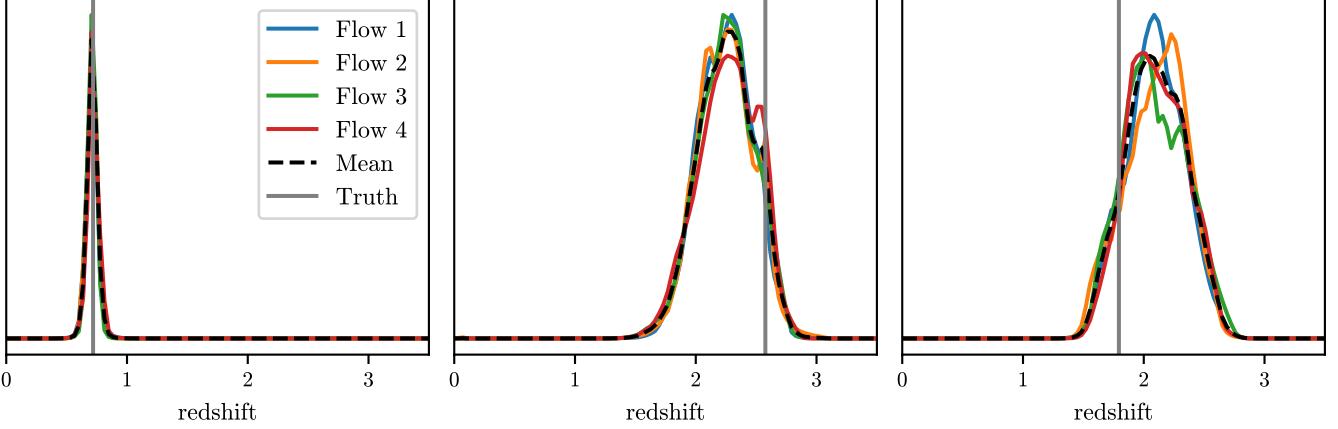


Figure 8. The ensemble of posteriors for three example galaxies. Flows 1-4 label the individual posteriors produced by each of the flows that make up the ensemble. The dashed black line is the mean of these individual posteriors and is the value used by the ensemble. The vertical gray line labeled “Truth” denotes the true redshift of the galaxy. Note these galaxies were specifically chosen for their broad, multimodal posteriors. The posteriors of most galaxies are sharp and unimodal.

chosen to model these additional two properties, any other values you desire can be similarly modeled.

For the latent distribution, we again use a Uniform distribution over the range $[-5, 5]$. For the bijection, we use

$$f = \text{RQ-RSC} \circ \text{Shift Bounds.} \quad (19)$$

The RQ-RSC acts on the two dimensional space of size and ellipticity, but also takes the galaxy redshift and magnitudes as inputs (see Figure 3). The redshifts and magnitudes are transformed to have zero mean and unit variance before being input to the neural network¹² that parameterizes the splines. Aside from the change in inputs, the RQ-RSC has the same settings as listed for the previous normalizing flow.

After training the flow (see Appendix A), we sample a size and ellipticity for each galaxy in the PZFlow catalog created in the previous section (conditioned on the true magnitudes), and plot the distribution of these features against the distribution in the test set (Figure 7). Once again, we see the normalizing flow does a good job of emulating the CosmoDC2 galaxy distribution.

The final simulated catalog consists of 10^4 galaxies, each with a redshift, $ugriz$ magnitudes including photometric errors, a true photo-z posterior consistent with the photometric errors, a size, and an ellipticity. This small catalog was generated for visualization purposes, but the normalizing flows can be used to generate catalogs of arbitrarily large size.

¹² These variables are standard scaled instead of mapped onto the domain $[-5, 5]$, because the neural network that parameterizes the splines has no limit on inputs, unlike the splines themselves, which are limited to the range $[-5, 5]$.

5. PHOTOMETRIC REDSHIFT ESTIMATION

In addition to forward modeling, normalizing flows are powerful and flexible models for density estimation. This makes them useful tools for estimating posterior distributions for galaxy properties, conditioned on observed features of the galaxy. A common example of this in cosmology is photometric redshift estimation, in which you estimate the redshift of a galaxy using its magnitude in several photometric bands. In this section, we use PZFlow as a photo-z estimator to demonstrate using normalizing flows for density estimation.

5.1. Training an Ensemble for photo-z estimation

When forward modeling in Section 4, we wanted a realistic model that captured the relevant correlations between galaxy photometry, redshift, shape, and size. However, when estimating redshifts, we do not simply want a realistic model, but rather a model that matches our specific galaxy sample as closely as possible.

When training deep learning models, the huge parameter space contains many different solutions, corresponding to different local minima in the parameter space. In the forward modeling application, we were content with finding a good local minimum, but in this application, we want to marginalize over the different potential models.

A full marginalization over the model parameters would be too computationally expensive, so instead we approximate this marginalization using an ensemble of normalizing flows. In other words, we train multiple normalizing flows under identical conditions, using different random initializations of the model parameters. This allows the optimization algorithm to explore different basins of attraction in the parameter space. In

the machine learning literature, this is known as a Deep Ensemble (Lakshminarayanan et al. 2017), and is a popular method for approximate bayesian marginalization (Wilson & Izmailov 2020; Fort et al. 2020).

We train an ensemble of 4 normalizing flows, each with the same architecture and training schedule as the regular flow described in Section 4. With PZFlow, this is as simple as swapping `FlowEnsemble` for `Flow` in the code.

For the training set, we use 100,000 galaxies from the catalog created in Section 4. Each galaxy in the training set has a true redshift and observed magnitudes in the *ugrizy* bands, with corresponding photometric errors. To account for the photometric error, at the start of each training epoch, we resample the training set from the photometric error distributions. In other words, each epoch, for each galaxy, we sample

$$m \sim p(\hat{m}, \sigma_m), \quad (20)$$

where \hat{m} are the observed magnitudes with photometric errors σ_m , and $p(\hat{m}, \sigma_m)$ is a Gaussian in flux space. This allows our ensemble of flows to approximate the distribution $p(z, m)$. For more details on training the ensemble, see Appendix A.

5.2. Estimating posteriors

For each flow in the ensemble, we estimate $p(z, m)$ by marginalizing over the photometric errors:

$$p(z|\hat{m}, \sigma_m) \propto \int dm p(z, m) p(m|\hat{m}, \sigma_m), \quad (21)$$

which is estimated by sampling $m \sim p(\hat{m}, \sigma_m)$ and averaging $p(z, m)$ over these samples. We then average the $p(z, m)$ from each flow, and normalize with respect to the redshift grid. This provides a redshift posterior for each galaxy.

Posteriors for three galaxies can be seen in Figure 8. Each flow produces a PDF which may contain slightly different features in each case. By averaging over the individual posteriors, we select for features that are common between models, while smoothing over features that are present in only a single model. We can also treat the ensemble of posteriors as a distribution over possible posteriors, which will allow for more consistent error calibration in cosmological analyses (Zhang et al. 2023).

We can calculate posteriors for galaxies with missing magnitudes by marginalizing over the missing magnitudes, as described in Section 4.2. An example is shown in Figure ??, which compares a posterior calculated using every LSST band to a posterior with the *u* band marginalized. Sacrificing the information in the *u* band increases the uncertainty of the redshift inference. In this instance, a second potential higher-redshift mode is

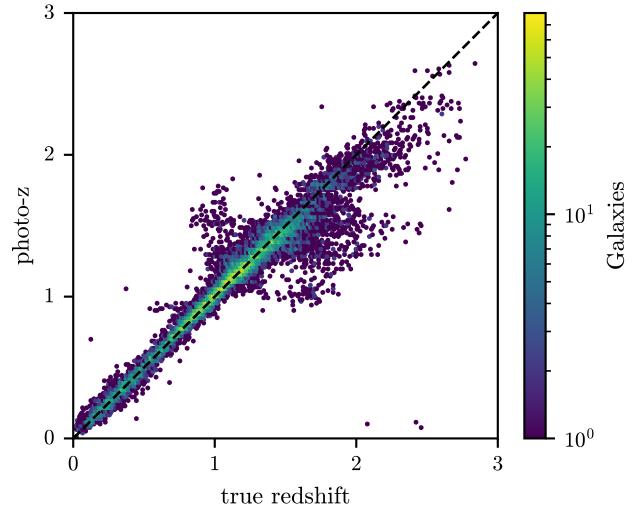


Figure 9. Photo-z point estimates (maximum a posteriori) vs true redshift for galaxies in the test set.

created. Conversely, this demonstrates how the addition of *u* band information rules out the higher redshift solution for this galaxy.

5.3. Photo-z metrics

In this section, we evaluate the performance of PZFlow using common photo-z metrics. Note these metrics are optimistic in the sense that the training set is representative of the test set, which is usually not the case in modern cosmology applications.

The most common metrics for photo-z estimation concern photo-z point estimates, which are a compression of the photo-z posterior to a single redshift estimate (e.g., Hildebrandt et al. 2010; Sánchez et al. 2014). We make the common choice of selecting the mode of the posteriors¹³. We compute metrics of the quantity $\Delta z = (z_{\text{phot}} - z_{\text{true}})/(1 + z_{\text{true}})$, where the denominator accounts for naturally greater uncertainties at high redshift.

Figure 9 compares the photo-z point estimates to the true redshifts. The point estimates for most galaxies lie along the diagonal, indicating strong performance. There are the common photo-z “wings”, indicating redshifts where important spectral features are transitioning between neighboring photometric bands. There is also a small population of high-redshift objects that were mistakenly identified as very low redshift objects. This point estimate plot is comparable to other high-performance machine learning photo-z esti-

¹³ The mean redshift is a poor choice, since photo-z posteriors are often multimodal, and so the mean value can lie between two modes at a redshift with very small probability density.

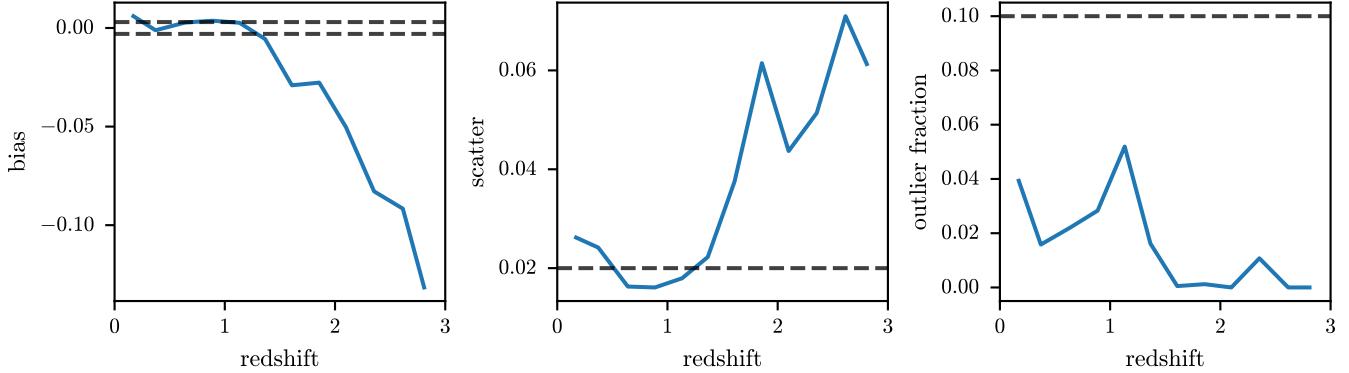


Figure 10. The bias, scatter, and outlier fraction of the photo-z point estimates as a function of true galaxy redshift. The dashed black lines represent the requirements for LSST cosmology as stated in the LSST DESC SRD (The LSST Dark Energy Science Collaboration et al. 2018). These lines are to provide a sense of scale for these metrics. You can see that PZFlow meets the bias and scatter requirements up to redshift ~ 1.5 , while meeting the outlier fraction requirements for all redshifts. We note that individual redshifts do not actually need to meet the bias requirement as long as the bias can be well calibrated via some other source, e.g. galaxy clustering.

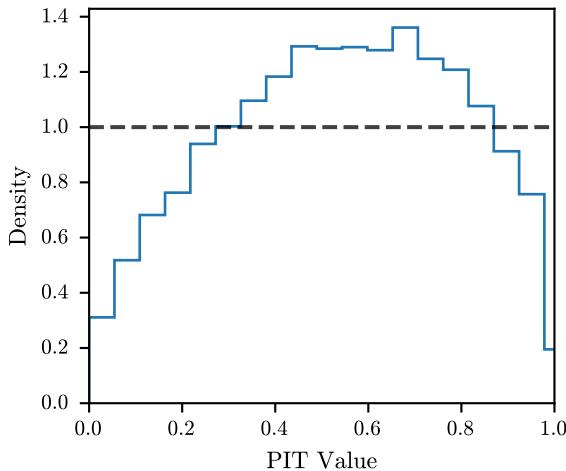


Figure 11. The probability integral transform (PIT) histogram for PZFlow photo-z posteriors. The PIT characterizes the calibration of the estimated posteriors, with the horizontal black line indicating perfect calibration.

mators when provided with representative training sets (Sánchez et al. 2014).

Figure 10 shows the photo-z point estimate metrics from the LSST DESC Science Requirements Document (SRD; The LSST Dark Energy Science Collaboration et al. 2018) as a function of true redshift. The *bias* is defined as the median of Δz ; the *scatter* is defined as $\text{IQR}/1.349$, where IQR is the interquartile range of Δz ; the *outlier fraction* is defined as the fraction of galaxies for which Δz is greater than three times the scatter. The requirements from the SRD are plotted in black to provide a sense of scale.

Like many photo-z estimators, PZFlow performs well to a redshift of approximately 1.5. At higher redshifts, our estimator does not meet the bias and scatter re-

quirements, because there is very little training data in this redshift range. We note however that for many cosmology applications, it is okay for the bias to exceed the required limits, as long as the bias can be well determined via some calibration process (Newman et al. 2015).

Another common metric is the probability integral transform (PIT) (see e.g. Schmidt & Malz et al. 2020, Dey et al. 2022). The PIT metric is used to determine if posteriors are well calibrated, i.e. if true values fall within X% confidence intervals X% of the time. The PIT histogram for our estimator is shown in Figure 11. Ideally, this histogram would be uniform and match the dashed horizontal line. The fact that the histogram bulges at the center indicates that our estimator is too conservative – i.e. the posteriors it produces are too broad. This can be explained by the fact that normalizing flows exhibit mode covering behavior (the opposite of the mode collapse seen in GANs; Salimans et al. 2016). In other words, because normalizing flows are trained by maximizing the likelihood of the training samples, they receive very high penalties for missing any modes in the data. As a result, they tend to conservatively spread out their density, in order to avoid missing any modes. This results in overly conservative posterior predictions. The low values at the edges of the PIT histogram indicate the relative rarity of catastrophic outliers, which is also reflected in the far right panel of Figure 10, where you can see that our estimator meets the requirement on the outlier fraction at all redshifts. There is also a slight rightward tilt. This indicates a small negative bias, which reflects the intrinsic prior towards smaller redshifts, as this is where the ma-

856 jority of galaxies in the training set lie. This negative
 857 bias is also visible in the far left panel of Figure 10.

858 The previous metrics analyze photo-z performance for
 859 point estimates, which are insufficient for modern cos-
 860 mology (Newman & Gruen 2022), and for ensembles of
 861 posteriors, which is often misleading and not a good in-
 862 dicator of performance for science applications (Schmidt
 863 & Malz et al. 2020). The methods introduced in this pa-
 864 per enable the creation of galaxy catalogs for which each
 865 galaxy has a *true* redshift posterior, which will enable
 866 more comprehensive evaluation of photo-z estimators.
 867 An full evaluation of photo-z estimators on a posterior-
 868 by-posterior basis is a major goal of the LSST DESC.

869 6. CONCLUSION & SUMMARY

870 In this paper we introduced PZFlow, a normalizing
 871 flow package in Python, designed for the statistical mod-
 872 eling of tabular astronomical data. We used PZFlow to
 873 forward model galaxy catalogs that include photometry,
 874 redshifts, sizes, and ellipticities. In addition, each galaxy
 875 in our catalog has a *true* redshift posterior, which can be
 876 convolved with measurement errors. These true posteri-
 877 ors allow a direct evaluation of the posteriors produced
 878 by photo-z estimators. A similar comparison can be
 879 made to posterior estimates for any other galaxy prop-
 880 erties.

881 Direct evaluation of photo-z posteriors is vital for fu-
 882 ture cosmology analyses which must use all of the infor-
 883 mation incorporated in the full redshift posteriors (New-
 884 man & Gruen 2022). This is because only the full red-
 885 shift posteriors allow one to account for degeneracies in
 886 color-redshift space, which will otherwise bias cosmo-
 887 logical inference. Previous evaluations of photo-z per-
 888 formance have focused on point estimates and metrics
 889 for ensembles of posteriors, but Schmidt & Malz et al.
 890 (2020) demonstrated these metrics are misleading and
 891 inadequate for modern cosmology. This work enables
 892 future analysis of photo-z estimators on a per-posterior
 893 basis, which is a major goal of the LSST DESC.

894 In addition to forward modeling galaxy catalogs, we
 895 demonstrated PZFlow’s utility as a density estimator
 896 that can be applied to photo-z estimation and other sta-
 897 tistical analyses. PZFlow achieves high accuracy with
 898 relatively little fine tuning and with very few modeling
 899 assumptions. This makes PZFlow a powerful tool for
 900 the statistical analysis of tabular astronomical data.

901 This paper was written using the showyourwork¹⁴
 902 workflow manager. The code to reproduce this paper
 903 is hosted publicly on Github¹⁵, and the code for each
 904 individual figure can be found by clicking on the Github
 905 logo in the margin next to that figure.

906 Thanks to Fran ois Lanusse for some early advice on
 907 normalizing flows and for reviewing the paper, and to
 908 Martine Lokken for testing PZFlow. J. F. Crenshaw,
 909 B. J. Kalmbach, and A. J. Connolly acknowledge sup-
 910 port from the DiRAC Institute in the Department of As-
 911 tronomy at the University of Washington. The DIRAC
 912 Institute is supported through generous gifts from the
 913 Charles and Lisa Simonyi Fund for Arts and Sciences,
 914 and the Washington Research Foundation. Z. Yan ac-
 915 knowledges support from the Max Planck Society and
 916 the Alexander von Humboldt Foundation in the frame-
 917 work of the Max Planck-Humboldt Research Award en-
 918 dowed by the German Federal Ministry of Education
 919 and Research. AIM acknowledges support during the
 920 course of this work from the Max Planck Society and the
 921 Alexander von Humboldt Foundation in the framework
 922 of the Max Planck-Humboldt Research Award endowed
 923 by the Federal Ministry of Education and Research.

924 Author contributions are listed below:

925 J. F. Crenshaw: created PZFlow, designed experiments,
 926 wrote paper.

927 J. B. Kalmbach: wrote early normalizing flow code that
 928 evolved into PZFlow; photo-z estimation.

929 A. Gagliano: validated code, developed use cases and
 930 associated tutorials; revised manuscript text.

931 Z. Yan: contributed to PZFlow and PhotErr.

932 A. J. Connolly: discussion during development.

933 A. I. Malz: consulted on design and testing of PZFlow.

934 S. J. Schimdt: consulted on design and testing of
 935 PZFlow and PhotErr.

936

937 *Software*: adam (Kingma & Ba 2015), corner
 938 (Foreman-Mackey 2016), dill (McKerns et al. 2012),
 939 jax (Bradbury et al. 2018), jupyter (Kluyver et al.
 940 2016), matplotlib (Hunter 2007), numpy (Harris et al.
 941 2020), pandas (Wes McKinney 2010; Reback et al. 2020),
 942 scipy (Virtanen et al. 2020), showyourwork (Luger et al.
 943 2021), scikit-learn (Pedregosa et al. 2011)

¹⁴ <https://show-your.work/>

¹⁵ <https://github.com/jfcrenshaw/pzflow-paper>

A. TRAINING DETAILS

In this section we list some technical details of training the normalizing flows. Every flow is trained via minimizing the negative log-likelihood

$$\mathcal{L} = -\mathbb{E}[\log p(x)], \quad (\text{A1})$$

where the expectation is performed over galaxies in the training set and $p(x)$ is defined in Equation 1.

For the main flow in Section 4, we trained for 150 epochs. We used the Adam optimizer (Kingma & Ba 2015), starting with a learning rate of 10^{-3} . We decreased the learning rate by a factor of 10 every 50 epochs. Training took 7 minutes on a Tesla P100 12GB GPU. The training loss for this flow is in the left panel of Figure 12.

For the conditional flow in 4, we trained for 450 epochs. Again, we used Adam with an initial learning rate of 10^{-3} . We decreased the learning rate by a factor of 10 every 150 epochs. The training loss for this flow is in the right panel of Figure 12.

For each of the flows that make up the flow ensemble in Section 5, we trained for 150 epochs using the Adam optimizer. We started each with a learning rate of 10^{-4} , and decreased the learning rate by a factor of 10 every 50 epochs. The training loss for the ensemble is in Figure 13. You can see that each of the flows achieves a different minimum loss. Apparently, each has found a different potential solution in the neural network’s parameter space.

B. LSST ERROR MODEL

We estimate photometric errors for LSST using a generalization of the error model from Ivezić et al. (2019). To derive the error model, we start with the noise-to-signal ratio (NSR) for an object with photon count C and background noise N_0 (which depends on seeing,

read-out noise, etc.):

$$\text{NSR}^2 = \frac{N_0^2 + C}{C^2}. \quad (\text{B2})$$

If we define $C = C_5$ when $\text{NSR} = 1/5$, then we can solve for N_0 and write

$$\text{NSR}^2 = \frac{1}{C_5} \left(\frac{C_5}{C} \right) + \left[\left(\frac{1}{5} \right)^2 - \frac{1}{C_5} \right] \left(\frac{C_5}{C} \right)^2. \quad (\text{B3})$$

Defining $x = C_5/C = 10^{(m-m_5)/2.5}$ and $\gamma = 1/5^2 - 1/C_5$, we have

$$\text{NSR}^2 = (0.04 - \gamma)x + \gamma x^2 \quad (\text{mag}^2), \quad (\text{B4})$$

which is Equation 5 from Ivezić et al. (2019). Values for the band-dependent parameter γ can be found in Table 2 of the same paper.

In the high signal-to-noise (SNR) limit, $\text{NSR} \ll 1$, and we can approximate

$$\sigma_{\text{rand}} = 2.5 \log_{10}(1 + \text{NSR}) \approx \text{NSR}. \quad (\text{B5})$$

This latter approximation is made by Ivezić et al. (2019), and errors are assumed to be Gaussian in magnitude space. In contrast, we use the exact form of Equation B5, and model errors as Gaussian in flux space. Note that after the photometric errors are applied, the error is re-calculated from the “observed” flux, and this new error is reported as the estimated photometric error. If the original photometric error were reported, it would provide a deterministic link to the original flux.

We have implemented this error model, along with several other extensions, in the Python package PhotErr, which is available on the Python Package Index¹⁶ (PyPI), and Github¹⁷. The extensions include different methods for handling non-detections, methods for modeling errors of extended objects (using models from van den Busch et al. 2020; Kuijken et al. 2019), and error models for the Roman and Euclid space telescopes (Spergel et al. 2015; Scaramella et al. 2022; Graham et al. 2020).

REFERENCES

- Bradbury, J., Frostig, R., Hawkins, P., et al. 2018, JAX: Composable Transformations of Python+NumPy Programs. <http://github.com/google/jax>
- Dacunha, T., Raveri, M., Park, M., Doux, C., & Jain, B. 2022, PhRvD, 105, 063529, doi: [10.1103/PhysRevD.105.063529](https://doi.org/10.1103/PhysRevD.105.063529)
- Dai, B., & Seljak, U. 2022, MNRAS, 516, 2363, doi: [10.1093/mnras/stac2010](https://doi.org/10.1093/mnras/stac2010)
- Dey, B., Zhao, D., Newman, J. A., et al. 2022, Calibrated Predictive Distributions via Diagnostics for Conditional Coverage, doi: [10.48550/arXiv.2205.14568](https://arxiv.org/abs/2205.14568)

¹⁶ <https://pypi.org/project/photerr/>

¹⁷ <https://github.com/jfcrenshaw/photerr>

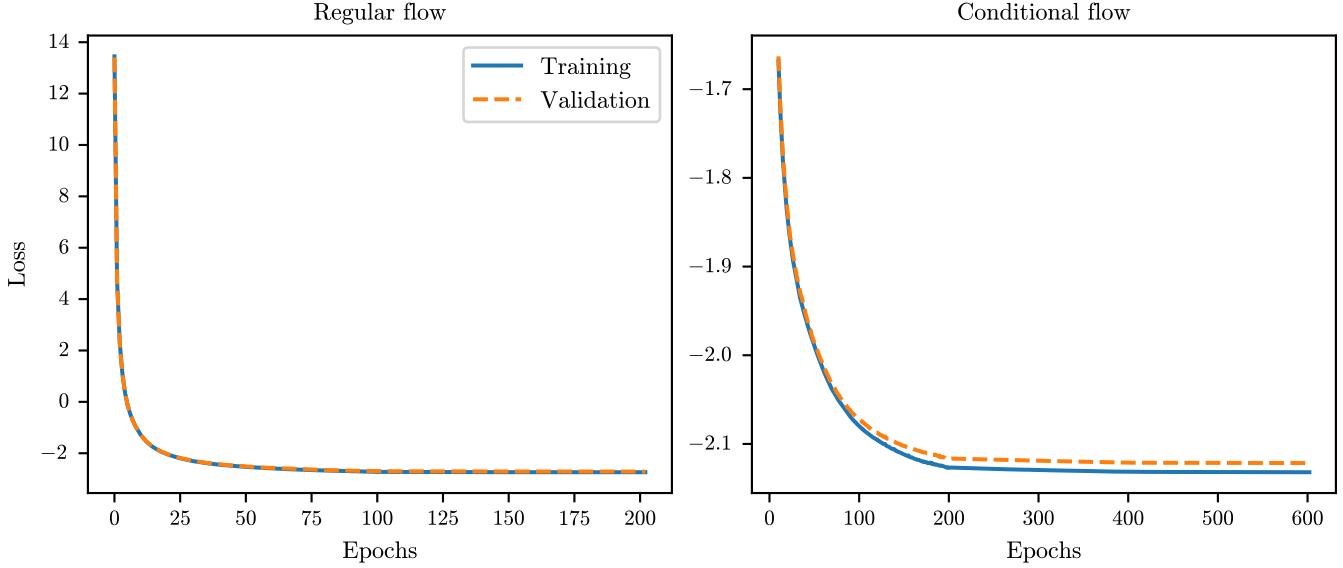


Figure 12. Training losses for the normalizing flows. Left: losses for the regular flow. After epochs 50 and 100, you can see a drop in the loss due to the decrease in the learning rate. Right: losses for the conditional flow. After epochs 150 and 300, you can see a drop in the loss due to the decrease in the learning rate.

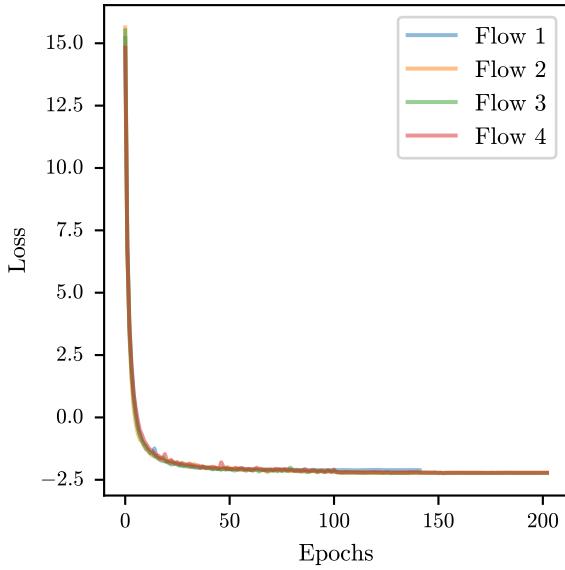


Figure 13. Training losses for the four flows in the flow ensemble. We have zoomed in to the bottom of the loss curve so you can see that each of the flows converges to a different minimum loss.

- 1036 Durkan, C., Bekasov, A., Murray, I., & Papamakarios, G.
 1037 2019, in Advances in Neural Information Processing
 1038 Systems 32, ed. H. M. Wallach, H. Larochelle,
 1039 A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, & R. Garnett
 1040 (Vancouver, Canada: Curran Associates, Inc.),
 1041 7511–7522. <https://arxiv.org/abs/1906.04032>
 1042 Falorsi, L., de Haan, P., Davidson, T. R., & Forré, P. 2019,
 1043 in Proceedings of Machine Learning Research, Vol. 89,
 1044 The 22nd International Conference on Artificial
 1045 Intelligence and Statistics, ed. K. Chaudhuri &
 1046 M. Sugiyama (Naha, Okinawa, Japan: PMLR),
 1047 3244–3253.
 1048 <http://proceedings.mlr.press/v89/falorsi19a.html>
 1049 Foreman-Mackey, D. 2016, The Journal of Open Source
 1050 Software, 1, 24, doi: [10.21105/joss.00024](https://doi.org/10.21105/joss.00024)
 1051 Fort, S., Hu, H., & Lakshminarayanan, B. 2020,
 1052 arXiv:1912.02757 [cs, stat].
 1053 <https://arxiv.org/abs/1912.02757>
 1054 Gemici, M. C., Rezende, D. J., & Mohamed, S. 2016,
 1055 CoRR, abs/1611.02304.
 1056 <https://arxiv.org/abs/1611.02304>
 1057 Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., et al. 2014,
 1058 CoRR, abs/1406.2661, doi: [10.48550/arXiv.1406.2661](https://doi.org/10.48550/arXiv.1406.2661)
 1059 Graham, M. L., Connolly, A. J., Ivezić, Ž., et al. 2018, The
 1060 Astronomical Journal, 155, 1,
 1061 doi: [10.3847/1538-3881/aa99d4](https://doi.org/10.3847/1538-3881/aa99d4)
 1062 Graham, M. L., Connolly, A. J., Wang, W., et al. 2020, The
 1063 Astronomical Journal, 159, 258,
 1064 doi: [10.3847/1538-3881/ab8a43](https://doi.org/10.3847/1538-3881/ab8a43)

- 1028 Dinh, L., Krueger, D., & Bengio, Y. 2015, in Proceedings of
 1029 the 3rd International Conference on Learning
 1030 Representations, ed. Y. Bengio & Y. LeCun, San Diego,
 1031 CA. <https://arxiv.org/abs/1410.8516>
 1032 Dinh, L., Sohl-Dickstein, J., & Bengio, S. 2017, in
 1033 Proceedings of the 5th International Conference on
 1034 Learning Representations, Toulon, France.
 1035 <https://arxiv.org/abs/1605.08803>

- 1065 Harris, C. R., Millman, K. J., van der Walt, S. J., et al.
 1066 2020, *Nature*, 585, 357, doi: [10.1038/s41586-020-2649-2](https://doi.org/10.1038/s41586-020-2649-2)
- 1067 Hassan, S., Villaescusa-Navarro, F., Wandelt, B., et al.
 1068 2022, *ApJ*, 937, 83, doi: [10.3847/1538-4357/ac8b09](https://doi.org/10.3847/1538-4357/ac8b09)
- 1069 Hildebrandt, H., Arnouts, S., Capak, P., et al. 2010, *A&A*,
 1070 523, A31, doi: [10.1051/0004-6361/201014885](https://doi.org/10.1051/0004-6361/201014885)
- 1071 Ho, J., Chen, X., Srinivas, A., Duan, Y., & Abbeel, P. 2019,
 1072 in Proceedings of Machine Learning Research, Vol. 97,
 1073 Proceedings of the 36th International Conference on
 1074 Machine Learning, ed. K. Chaudhuri & R. Salakhutdinov
 1075 (Long Beach, CA, USA: PMLR), 2722–2730.
 1076 <http://proceedings.mlr.press/v97/ho19a.html>
- 1077 Hoogeboom, E., Cohen, T. S., & Tomczak, J. M. 2020,
 1078 CoRR, abs/2001.11235.
<https://arxiv.org/abs/2001.11235>
- 1080 Hunter, J. D. 2007, *Computing in Science & Engineering*, 9,
 1081 90, doi: [10.1109/MCSE.2007.55](https://doi.org/10.1109/MCSE.2007.55)
- 1082 Ivezić, Ž., Kahn, S. M., Tyson, J. A., et al. 2019, *The
 1083 Astrophysical Journal*, 873, 111,
 1084 doi: [10.3847/1538-4357/ab042c](https://doi.org/10.3847/1538-4357/ab042c)
- 1085 Jaini, P., Kobyzhev, I., Yu, Y., & Brubaker, M. 2020, in
 1086 Proceedings of Machine Learning Research, Vol. 119,
 1087 arXiv:1907.04481 [Cs, Math, Stat] (Virtual: PMLR),
 1088 4673–4681. <https://arxiv.org/abs/1907.04481>
- 1089 Kessler, R., Narayan, G., Avelino, A., et al. 2019,
 1090 *Publications of the Astronomical Society of the Pacific*,
 1091 131, 094501, doi: [10.1088/1538-3873/ab26f1](https://doi.org/10.1088/1538-3873/ab26f1)
- 1092 Kingma, D. P., & Ba, J. 2015, in 3rd International
 1093 Conference on Learning Representations, ed. Y. Bengio &
 1094 Y. LeCun, San Diego, CA.
<https://arxiv.org/abs/1412.6980>
- 1096 Kingma, D. P., & Welling, M. 2014, in 2nd International
 1097 Conference on Learning Representations, ed. Y. Bengio &
 1098 Y. LeCun, Banff, AB, Canada,
 1099 doi: [10.48550/arXiv.1312.6114](https://doi.org/10.48550/arXiv.1312.6114)
- 1100 Kluyver, T., Ragan-Kelley, B., Pérez, F., et al. 2016, in
 1101 Positioning and Power in Academic Publishing: Players,
 1102 Agents and Agendas, ed. F. Loizides & B. Scmidt
 1103 (Netherlands: IOS Press), 87–90.
<https://eprints.soton.ac.uk/403913/>
- 1105 Kobyzhev, I., Prince, S. J. D., & Brubaker, M. A. 2020,
 1106 *IEEE Trans. Pattern Anal. Mach. Intell.*, 1,
 1107 doi: [10.1109/TPAMI.2020.2992934](https://doi.org/10.1109/TPAMI.2020.2992934)
- 1108 Korytov, D., Hearin, A., Kovacs, E., et al. 2019, *The
 1109 Astrophysical Journal Supplement Series*, 245, 26,
 1110 doi: [10.3847/1538-4365/ab510c](https://doi.org/10.3847/1538-4365/ab510c)
- 1111 Kuijken, K., Heymans, C., Dvornik, A., et al. 2019, *A&A*,
 1112 625, A2, doi: [10.1051/0004-6361/201834918](https://doi.org/10.1051/0004-6361/201834918)
- 1113 Lakshminarayanan, B., Pritzel, A., & Blundell, C. 2017, in
 1114 Advances in Neural Information Processing Systems 30:
 1115 Annual Conference on Neural Information Processing
 1116 Systems 2017, December 4–9, 2017, Long Beach, CA,
 1117 USA, ed. I. Guyon, U. von Luxburg, S. Bengio, H. M.
 1118 Wallach, R. Fergus, S. V. N. Vishwanathan, &
 1119 R. Garnett, 6402–6413.
<https://proceedings.neurips.cc/paper/2017/hash/9ef2ed4b7fd2c810847ffa5fa85bce38-Abstract.html>
- 1120 Lokken, M., Gagliano, A., Narayan, G., et al. 2022, *The
 1121 Simulated Catalogue of Optical Transients and
 1122 Correlated Hosts (SCOTCH)*.
<https://ui.adsabs.harvard.edu/abs/2022arXiv220602815L>
- 1123 LSST Dark Energy Science Collaboration, Abolfathi, B.,
 1124 Alonso, D., et al. 2021, *The Astrophysical Journal
 1125 Supplement Series*, 253, 31,
 1126 doi: [10.3847/1538-4365/abd62c](https://doi.org/10.3847/1538-4365/abd62c)
- 1127 Luger, R., Bedell, M., Foreman-Mackey, D., et al. 2021,
 1128 Mapping Stellar Surfaces III: An Efficient, Scalable, and
 1129 Open-Source Doppler Imaging Model.
<https://ui.adsabs.harvard.edu/abs/2021arXiv211006271L>
- 1130 Malz, A. I., Lanusse, F., Crenshaw, J. F., & Graham, M. L.
 1131 2021, An Information-Based Metric for Observing
 1132 Strategy Optimization, Demonstrated in the Context of
 1133 Photometric Redshifts with Applications to Cosmology,
 1134 doi: [10.48550/arXiv.2104.08229](https://doi.org/10.48550/arXiv.2104.08229)
- 1135 Mandelbaum, R., Seljak, U., Hirata, C. M., et al. 2008,
 1136 Monthly Notices of the Royal Astronomical Society, 386,
 1137 781, doi: [10.1111/j.1365-2966.2008.12947.x](https://doi.org/10.1111/j.1365-2966.2008.12947.x)
- 1138 McKerns, M. M., Strand, L., Sullivan, T., Fang, A., &
 1139 Aivazis, M. A. G. 2012, CoRR, abs/1202.1056.
<http://arxiv.org/abs/1202.1056>
- 1140 Newman, J. A., & Gruen, D. 2022, *Annu. Rev. Astron.
 1141 Astrophys.*, 60, annurev,
 1142 doi: [10.1146/annurev-astro-032122-014611](https://doi.org/10.1146/annurev-astro-032122-014611)
- 1143 Newman, J. A., Abate, A., Abdalla, F. B., et al. 2015,
 1144 *Astroparticle Physics*, 63, 81,
 1145 doi: [10.1016/j.astropartphys.2014.06.007](https://doi.org/10.1016/j.astropartphys.2014.06.007)
- 1146 Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011,
 1147 *Journal of Machine Learning Research*, 12, 2825.
<http://jmlr.org/papers/v12/pedregosa11a.html>
- 1148 Reback, J., McKinney, W., jbrockmendel, et al. 2020,
 1149 Pandas-Dev/Pandas: Pandas 1.0.3, Zenodo.
<https://doi.org/10.5281/zenodo.3715232>
- 1150 Rezende, D. J., Papamakarios, G., Racanière, S., et al.
 1151 2020, in Proceedings of Machine Learning Research, Vol.
 1152 119, arXiv:2002.02428 [Cs, Stat] (Virtual: PMLR),
 1153 8083–8092. <https://arxiv.org/abs/2002.02428>

- 1161 Salimans, T., Goodfellow, I. J., Zaremba, W., et al. 2016, in
 1162 Advances in Neural Information Processing Systems 29:
 1163 Annual Conference on Neural Information Processing
 1164 Systems 2016, December 5-10, 2016, Barcelona, Spain,
 1165 ed. D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon,
 1166 & R. Garnett, 2226–2234.
 1167 <https://proceedings.neurips.cc/paper/2016/hash/8a3363abe792db2d8761d6403605aeb7-Abstract.html>
- 1168 Sánchez, C., Carrasco Kind, M., Lin, H., et al. 2014,
 1170 Monthly Notices of the Royal Astronomical Society, 445,
 1171 1482, doi: [10.1093/mnras/stu1836](https://doi.org/10.1093/mnras/stu1836)
- 1172 Scaramella, R., Amiaux, J., Mellier, Y., et al. 2022, A&A,
 1173 662, A112, doi: [10.1051/0004-6361/202141938](https://doi.org/10.1051/0004-6361/202141938)
- 1174 Schmidt, S. J., Malz, A. I., Soo, J. Y. H., et al. 2020,
 1175 arXiv:2001.03621 [astro-ph].
 1176 <https://arxiv.org/abs/2001.03621>
- 1177 Spergel, D., Gehrels, N., Baltay, C., et al. 2015, Wide-Field
 1178 Infrared Survey Telescope-Astrophysics Focused
 1179 Telescope Assets WFIRST-AFTA 2015 Report,
 1180 doi: [10.48550/arXiv.1503.03757](https://arxiv.org/abs/1503.03757)
- 1181 Stylianou, N., Malz, A. I., Hatfield, P., Crenshaw, J. F., &
 1182 Gschwend, J. 2022, Publications of the Astronomical
 1183 Society of the Pacific, 134, 044501,
 1184 doi: [10.1088/1538-3873/ac59bf](https://doi.org/10.1088/1538-3873/ac59bf)
- 1185 The LSST Dark Energy Science Collaboration,
 1186 Mandelbaum, R., Eifler, T., et al. 2018, arXiv e-prints,
 1187 1809, arXiv:1809.01669.
- 1188 <http://adsabs.harvard.edu/abs/2018arXiv180901669T>
- 1189 van den Busch, J. L., Hildebrandt, H., Wright, A. H., et al.
 1190 2020, Astronomy and Astrophysics, 642, A200,
 1191 doi: [10.1051/0004-6361/202038835](https://doi.org/10.1051/0004-6361/202038835)
- 1192 Virtanen, P., Gommers, R., Oliphant, T. E., et al. 2020,
 1193 Nature Methods, 17, 261, doi: [10.1038/s41592-019-0686-2](https://doi.org/10.1038/s41592-019-0686-2)
- 1194 Wes McKinney. 2010, in Proceedings of the 9th Python in
 1195 Science Conference, ed. S. van der Walt & Jarrod
 1196 Millman, 56–61, doi: [10.25080/Majora-92bf1922-00a](https://doi.org/10.25080/Majora-92bf1922-00a)
- 1197 Wilson, A. G., & Izmailov, P. 2020, arXiv:2002.08791 [cs,
 1198 stat]. <https://arxiv.org/abs/2002.08791>
- 1199 Winkler, C., Worrall, D., Hoogeboom, E., & Welling, M.
 1200 2019, CoRR, abs/1912.00042.
 1201 <https://arxiv.org/abs/1912.00042>
- 1202 Zhang, T., Rau, M. M., Mandelbaum, R., Li, X., & Moews,
 1203 B. 2023, Monthly Notices of the Royal Astronomical
 1204 Society, 518, 709, doi: [10.1093/mnras/stac3090](https://doi.org/10.1093/mnras/stac3090)