

---

# Squared Neural Families: A New Class of Tractable Density Models

---

**Russell Tsuchida\***  
Data61-CSIRO

**Cheng Soon Ong**  
Data61-CSIRO &  
Australian National University

**Dino Sejdinovic**  
School of CMS & AIML,  
The University of Adelaide

## Abstract

Flexible models for probability distributions are an essential ingredient in many machine learning tasks. We develop and investigate a new class of probability distributions, which we call a Squared Neural Family (SNEFY), formed by squaring the 2-norm of a neural network and normalising it with respect to a base measure. Following the reasoning similar to the well established connections between infinitely wide neural networks and Gaussian processes, we show that SNEFYs admit closed form normalising constants in many cases of interest, thereby resulting in flexible yet fully tractable density models. SNEFYs strictly generalise classical exponential families, are closed under conditioning, and have tractable marginal distributions. Their utility is illustrated on a variety of density estimation, conditional density estimation, and density estimation with missing data tasks.

## 1 Introduction

Probabilistic modelling lies at the heart of machine learning. In both traditional and contemporary settings, ensuring the probability model is appropriately normalised (or otherwise bypassing the need to compute normalising constants) is of central interest for maximum likelihood estimation and related statistical inference procedures. Tractable normalising constants allow the use of probability models in a variety of applications, such as anomaly detection, denoising and generative modelling.

Traditional statistical approaches [27, 5] rely on mathematically convenient models such as exponential families [50]. Such models can often be normalised in closed form, but are often only suitable for relatively low-dimensional and simple data. Early approaches in deep learning use neural networks as energy functions inside Gibbs distributions [21, Equation 2]. Such distributions typically have very intractable normalising constants, and so either surrogate losses for the negative log likelihood involving the energy function are used or MCMC, score matching [16], variational or sampling methods [15, 30] are used to approximate the normalising constant. See [35, §2.6.1, §2.2] for more energy-based and other methods.

Modern computer technology and power allows us more flexible models [8, page 68], by partially or completely relaxing the requirement for mathematical tractability. Contemporary high-dimensional modelling, such as generative image models, rely primarily on neural network models [18, 11, 35, 34]. For example, normalising flows [40] use constrained neural network layers to transform random variables from base measures with tractable Jacobian determinants associated with appropriately constrained but flexible pushforward maps. Similarly to normalising flows, we use neural networks to define expressive densities, but we model the densities directly without using constrained transformations of random variables. Moreover, our approach can be readily applied in conjunction with normalising flows and other deep learning devices such as deep feature extractors.

---

\*Software available at <https://github.com/RussellTsuchida/snefy>.

Support $\mathbb{X}$	Base measure $\mu$	Sufficient statistic $t(\mathbf{x})$	Activation function $\sigma$	$\mathbf{b} \neq \mathbf{0}$	Kernel $k_{\sigma,t,\mu}$
<b>Any setting mirroring a previously derived NNGPK, e.g.</b>					
$\mathbb{R}^d$	$\Phi_{\mathbf{C},\mathbf{0}}$	$\text{Id}(\mathbf{x})$	erf	$\times$	[52]
			$(\cdot)_+^p, p \in \mathbb{N}$	$\times$	[4]
			LRReLU	$\times$	[48]
			GELU [14]	$\times$	[47]
			cos	$\times$	[37]
<b>Any tractable exponential family [50, 32] setting, e.g.</b>					
$\mathbb{R}^d$	$\Phi_{\mathbf{C},\mathbf{m}}$	$\frac{\mathbf{x}}{\ \mathbf{x}\ _2}$	exp	$\checkmark$	Kernel 3
$\mathbb{S}^{d-1}$	Uniform	Id( $\mathbf{x}$ )		$\checkmark$	
$\mathbb{R}^d$	$\Phi_{\mathbf{C},\mathbf{m}}$			$\checkmark$	Kernel 7
$\{0, 1, 2, \dots\}$	$(x!)^{-1}\nu$	$\checkmark$		Kernel 8	
<b>New tractable integration settings, e.g.</b>					
$\mathbb{R}^d$	$\Phi_{\mathbf{C},\mathbf{m}}$	Id( $\mathbf{x}$ )	cos	$\checkmark$	Kernel 2
$[0, 1]^d$	Uniform	$\Phi^{-1}(\mathbf{x})$			
$\mathbb{R}^d$	$\Phi_{\mathbf{C},\mathbf{m}}$	Id( $\mathbf{x}$ )	Snake <sub>a</sub> [55, 39]	$\checkmark$	Kernel 6

Table 1: Examples of settings admitting a closed-form for the normalising constant  $z(\mathbf{V}, \Theta)$  (2) by leveraging a closed-form NNK  $k_{\sigma,t,\mu}$  (4). In each case,  $z(\mathbf{V}, \Theta) = \text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)$ , where the entries of the matrix  $\mathbf{K}_\Theta$  are described according to the NNK  $k_{\sigma,t,\mu}$  in Identity 1.  $\Phi_{\mathbf{C},\mathbf{m}}$  denotes the CDF of a multivariate normal distribution with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$  and  $\nu$  denotes counting measure. Rows with citations have been considered previously in the context of NNGPKs, but not as normalising constants and not with a reversal of the role of input and parameter. Note the cases where  $\mathbf{b} \neq \mathbf{0}$ ; this setting is not considered by others, because when the role of parameters and data is in the usual setting,  $\mathbf{b} = \mathbf{0}$  covers a sufficiently general setting. Noticing that SNEFYs strictly generalise exponential family mixture models (see § 3.2), fixing  $\sigma(\cdot) = \exp(\cdot/2)$  and a given base measure  $\mu$  and sufficient statistic  $t$  for which the exponential family log-partition function is known also leads to tractable normalising constants. More known closed-form kernels as well as approximate kernels that can be adapted to our setting are given in [12].

**Contributions.** Let  $(\Omega, \mathcal{F}, \mu)$  denote a measure space with  $\Omega \subseteq \mathbb{R}^d$ , sigma algebra  $\mathcal{F}$ , and nonnegative measure  $\mu$ . Let  $\mathbf{V} \in \mathbb{R}^{m \times n}$  be the readout parameters of a neural network and let  $\mathbf{W} \in \mathbb{R}^{n \times D}$  and  $\mathbf{b} \in \mathbb{R}^n$  be the weights and biases of a hidden layer of a neural network  $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^m$  with activation function  $\sigma$ . For some support  $\mathbb{X} \in \mathcal{F}$ , define a probability distribution  $P$  (and corresponding probability density  $p$  with respect to base measure  $\mu$ ) to be proportional to squared 2-norm of the evaluation of the neural network  $\mathbf{f}$ ,

$$P(d\mathbf{x}; \mathbf{V}, \Theta) \triangleq \frac{\mu(d\mathbf{x})}{z(\mathbf{V}, \Theta)} \|\mathbf{f}(t(\mathbf{x}); \mathbf{V}, \Theta)\|_2^2, \quad \mathbf{f}(t; \mathbf{V}, \Theta) = \mathbf{V}\sigma(\mathbf{W}t + \mathbf{b}), \quad \Theta = (\mathbf{W}, \mathbf{b}), \quad (1)$$

whenever the normalising constant  $z(\mathbf{V}, \Theta) \triangleq \int_{\mathbb{X}} \|\mathbf{f}(t(\mathbf{x}); \mathbf{V}, \Theta)\|_2^2 \mu(d\mathbf{x})$  is finite and non-zero. Here we call  $\mu$  the base measure,  $t : \mathbb{X} \rightarrow \mathbb{R}^D$  the sufficient statistic<sup>2</sup> and  $\sigma$  the activation function. We will call the corresponding family of probability distributions, parametrised by  $(\mathbf{V}, \Theta)$ , a *squared neural family* (SNEFY) on  $\mathbb{X}$ , and denote it by  $\text{SNEFY}_{\mathbb{X},t,\sigma,\mu}$ . SNEFYs are new flexible probability distribution models that strike a good balance between mathematical tractability, expressivity and computational efficiency. When a random vector  $\mathbf{x}$  follows a SNEFY distribution indexed by parameters  $(\mathbf{V}, \Theta)$ , we write  $\mathbf{x} \sim \text{SNEFY}_{\mathbb{X},t,\sigma,\mu}(\mathbf{V}, \Theta)$  or simply  $\mathbf{x} \sim P(\cdot; \mathbf{V}, \Theta)$ , where there is no ambiguity. Our main technical challenge is in exactly computing the normalising constant,  $z(\mathbf{V}, \Theta)$ , where

$$z(\mathbf{V}, \Theta) \triangleq \int_{\mathbb{X}} \|\mathbf{f}(t(\mathbf{x}); \mathbf{V}, \Theta)\|_2^2 \mu(d\mathbf{x}), \quad \mathbf{f}(t; \mathbf{V}, \Theta) = \mathbf{V}\sigma(\mathbf{W}t + \mathbf{b}), \quad \Theta = (\mathbf{W}, \mathbf{b}). \quad (2)$$

The normalising constants we consider are special cases in the sense that they apply to specific (but commonly appearing in applications) choices of activation function  $\sigma$ , sufficient statistic  $t$  and base measure  $\mu$  over support  $\mathbb{X}$ . See Table 1. Our analysis both exploits and informs a connection with

<sup>2</sup>We later verify that  $t$  is indeed a sufficient statistic, see (7). Note  $t$  maps from  $d$  to  $D$  dimensions.

so-called neural network Gaussian process kernels (NNGPKs) [31] in a generalised form, which we refer to as neural network kernels (NNKs).

We discuss some important theoretical properties of SNEFY such as exact normalising constant calculation, marginal distributions, conditional distributions, and connections with other probability models. We then consider a deep learning setting, where SNEFYs can either be used as base distributions in (non-volume preserving) normalising flows [46], or may describe flexible conditional density models with deep learning feature extractors. We demonstrate SNEFY on a variety of datasets.

## 1.1 Background

**Notation** We use lower case non-bold (like  $a$ ) to denote scalars, lower case bold to denote vectors (like  $\mathbf{a}$ ) and upper case bold to denote matrices (like  $\mathbf{A}$ ). Random variables (scalar or vector) are additionally typeset in sans-serif (like  $a$  and  $\mathbf{a}$ ). The special zero vector  $(0, \dots, 0)$  and identity matrix elements are  $\mathbf{0}$  and  $\mathbf{I}$ . We use subscripts to extract (groups of) indices, so that for example,  $w_i$  is the  $i$ th row of the matrix  $\mathbf{W}$  (as a column vector),  $v_{\cdot,i}$  is the  $i$ th column of the matrix  $\mathbf{V}$ , and  $b_i$  is the  $i$ th element of the vector  $\mathbf{b}$ . We use  $\Theta = (\mathbf{W}, \mathbf{b}) \in \mathbb{R}^{n \times (D+1)}$  to denote the concatenated hidden layer weights and biases. Correspondingly, we write  $\theta_i = (w_i, b_i) \in \mathbb{R}^{D+1}$  for the  $i$ th row of  $\Theta$ . We will use a number of special functions.  $\Phi_{\mathbf{C}, \mathbf{m}}$  denotes the cumulative distribution function (CDF) of a Gaussian random vector with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$ . We also use a shorthand  $\Phi_{\mathbf{C}} = \Phi_{\mathbf{C}, \mathbf{0}}$  and  $\Phi = \Phi_{\mathbf{I}, \mathbf{0}}$ .

**Single hidden layer neural networks** We consider a feedforward neural network  $f: \mathbb{R}^D \rightarrow \mathbb{R}^m$ ,

$$f(\mathbf{t}; \mathbf{V}, \Theta) = \mathbf{V} \sigma(\mathbf{W} \mathbf{t} + \mathbf{b}) \quad (3)$$

with activation function  $\sigma$ , hidden weights  $\mathbf{W} \in \mathbb{R}^{n \times D}$  and biases  $\mathbf{b} \in \mathbb{R}^n$  and readout parameters  $\mathbf{V} \in \mathbb{R}^{m \times n}$ . Here  $\sigma$  is applied element-wise to its vector inputs, returning a vector of the same shape.

**Neural network kernels** In certain theories of deep learning, one often encounters a bivariate Gaussian integral called the *neural network Gaussian process kernel* NNGPK. The NNGPK first arose as the covariance function of a well-behaved single layer neural network with random weights [31]. In the limit as the width of the network grows to infinity, the neural network (3) with suitably well-behaved (say, independent Gaussian) random weights converges to a zero-mean Gaussian process, so that the NNGPK characterises the law of the neural network predictions. These limiting models can be used as functional priors in a classical Gaussian process sense [53].

In our setting, the positive semidefinite (PSD) NNGPK appears in an entirely novel context, where the role of the hidden weights and biases  $\theta_i = (w_i, b_i)$  and the data  $\mathbf{x}$  is reversed. Instead of marginalising out the parameters and evaluating at the data, we marginalise out the data and evaluate at the parameters. The NNGPK  $k_{\sigma, \text{Id}, \Phi_{\mathbf{I}}}$  admits a representation of the form

$$k_{\sigma, \text{Id}, \Phi_{\mathbf{I}}}(\theta_i, \theta_j) \triangleq \mathbb{E}_{\mathbf{x}}[\sigma(\mathbf{w}_i^\top \mathbf{x} + b_i) \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)], \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (4)$$

We do not discuss in detail how it is constructed in earlier works [31, 22, 33, 17], where usually  $b_i = b_j = 0^3$ , but not always [49]. When  $b_i = b_j = 0$ , closed-form expressions for the NNGPK are available for different choices of  $\sigma$  and  $\mathbf{t}$  [52, 20, 4, 48, 37, 47, 28, 13]. However, the setting of  $\mathbf{b} \neq \mathbf{0}$  is important in our context (as we show in §2.2) and presents additional analytical challenges.

We will require a more general notion of an NNGPK which we call a *neural network kernel* (NNK). We introduce a function  $\mathbf{t}$  which may be thought of as a warping function applied to the input data. Such warping is common in kernels and covariance functions and can be used to induce desirable analytical and practical properties [37, 29, 24, §5.4.3]. We also integrate with respect to more general measures  $\mu$  instead of the standard Gaussian CDF,  $\Phi$ . We define the NNK to be

$$k_{\sigma, \mathbf{t}, \mu}(\theta_i, \theta_j) \triangleq \int_{\mathbb{X}} \sigma(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{w}_j^\top \mathbf{t}(\mathbf{x}) + b_j) \mu(d\mathbf{x}). \quad (5)$$

<sup>3</sup>After accounting for the reversal of parameters and data.

## 2 Closed form squared neural families

### 2.1 Normalising constants

Observe from (2) that by swapping the order of integration and multiplication by  $\mathbf{V}$ , the normalising constant is quadratic in elements of  $\mathbf{V}$ . The coefficients of the quadratic depend on  $\Theta = (\mathbf{W}, \mathbf{b})$ . We now characterise these coefficients of the quadratic in terms of the NNK evaluated at rows  $\theta_i, \theta_j$  of  $\Theta$ , which are totally independent of  $\mathbf{V}$ . The proof of the following is given in Appendix A.

**Identity 1.** *The integral (2) admits a representation of the form*

$$z(\mathbf{V}, \Theta) = \text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta) \quad (6)$$

where  $k_{\sigma, \mathbf{t}, \mu}$  is as defined in (5), and  $\mathbf{K}_\Theta$  is the PSD matrix whose  $ij$ th entry is  $k_{\sigma, \mathbf{t}, \mu}(\theta_i, \theta_j)$ .

By Identity 1, the normalised measure (1) then admits the explicit representation

$$P(d\mathbf{x}; \mathbf{V}, \Theta) = \frac{\text{Tr}(\mathbf{V}^\top \mathbf{V} \widetilde{\mathbf{K}}_\Theta(\mathbf{x}))}{\text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)} \mu(d\mathbf{x}) = \frac{\text{vec}(\mathbf{V}^\top \mathbf{V})^\top \text{vec}(\widetilde{\mathbf{K}}_\Theta(\mathbf{x}))}{\text{vec}(\mathbf{V}^\top \mathbf{V})^\top \text{vec}(\mathbf{K}_\Theta)} \mu(d\mathbf{x}),$$

where  $\widetilde{\mathbf{K}}_\Theta(\mathbf{x})$  is the PSD matrix whose  $ij$ th entry is  $\sigma(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{w}_j^\top \mathbf{t}(\mathbf{x}) + b_j)^\top$ . We used the cyclic property of the trace, writing the numerator as  $\text{Tr}(\boldsymbol{\sigma}^\top \mathbf{V}^\top \mathbf{V} \boldsymbol{\sigma}) = \text{Tr}(\mathbf{V} \boldsymbol{\sigma} \boldsymbol{\sigma}^\top \mathbf{V}^\top) = \text{Tr}(\mathbf{V}^\top \mathbf{V} \boldsymbol{\sigma} \boldsymbol{\sigma}^\top)$ . We emphasise again that the role of the data  $\mathbf{t}(\mathbf{x})$  and the hidden weights and biases  $\Theta$  in the NNK  $k_{\sigma, \mathbf{t}, \mu}$  are reversed compared with how they have appeared in previous settings. We may compute evaluations of  $k_{\sigma, \mathbf{t}, \mu}$  in closed form for special cases of  $(\sigma, \mathbf{t}, \mu)$  using various identities in  $\mathcal{O}(d)$ , where  $d$  is the dimensionality of the domain of integration, as we soon detail in § 2.2. Combined with the trace inner product, this leads to a total cost of computing  $z(\mathbf{V}, \Theta)$  of  $\mathcal{O}(m^2 n + dn^2)$ , where  $n$  and  $m$  are respectively the number of neurons in the first and second layers.

**Remark 1** (Alternative parameterisations). *SNEFY models depend on readout parameters  $\mathbf{V}$  only through the direction of  $\text{vec}(\mathbf{V}^\top \mathbf{V})$  and not on its norm or sign. For example, one can always find another parameterisation of readout parameters that results in the same probability distribution but has a normalising constant of 1. Furthermore, noticing that  $\mathbf{V}$  only appears as a PSD matrix  $\mathbf{M} \triangleq \mathbf{V}^\top \mathbf{V}$  of rank at most  $\min(m, n)$ , one may alternatively parameterise a SNEFY by  $(\mathbf{M}, \Theta)$ .*

### 2.2 Neural network kernels

In § 2.1, we reduced computation of the integral (2) to computation of a quadratic involving evaluations of the NNK (5). Several closed-forms are known for different settings of  $\sigma$ , all with  $\mathbf{t} = \text{Id}$  and  $\mathbf{b} = \mathbf{0}$ . The motivation behind derivation of existing known results is from the perspective of inference in infinitely wide Bayesian neural networks [31] or to derive certain integrals involved in computing predictors of infinitely wide neural networks trained using gradient descent [17]. Here we describe some new settings that have not been investigated previously that are useful for the new setting of SNEFY. Recall from (5), that our kernel,  $k_{\sigma, \mathbf{t}, \mu}$ , is parametrised by activation function  $\sigma$ , warping function (sufficient statistic)  $\mathbf{t}$ , and base measure  $\mu$ . All derivations are given in Appendix B.

The first kernel describes how we may express the kernels with arbitrary Gaussian base measures  $\Phi_{\mathbf{C}, m}$  in terms of the kernels with isotropic Gaussian base measures  $\Phi$ . This means that it suffices to consider isotropic Gaussian base measures in place of arbitrary Gaussian base measures.

**Kernel 1.**  $k_{\sigma, \text{Id}, \Phi_{\mathbf{C}, m}}(\theta_i, \theta_j) = k_{\sigma, \text{Id}, \Phi}(\mathcal{T}\theta_i, \mathcal{T}\theta_j)$ , where  $\mathcal{T}\Theta = (\mathbf{W}\mathbf{A}, \mathbf{b} + \mathbf{W}\mathbf{m})$ ,  $\mathcal{T}\theta_i = (\mathbf{w}_i^\top \mathbf{A}, b_i + \mathbf{w}_i^\top \mathbf{m})$  and  $\mathbf{A}$  is a matrix factor such that covariance  $\mathbf{C} = \mathbf{A}\mathbf{A}^\top$ .

This kernel can also be used to describe kernels corresponding with Gaussian mixture model base measures. The second kernel we describe is a minor extension to the case  $\mathbf{b} \neq \mathbf{0}$  of a previously considered kernel [37].

**Kernel 2.**  $k_{\cos, \text{Id}, \Phi}(\theta_i, \theta_j) = \frac{\cos |b_i - b_j|}{2} \exp\left(\frac{-\|\mathbf{w}_j - \mathbf{w}_i\|^2}{2}\right) + \frac{\cos |b_i + b_j|}{2} \exp\left(\frac{-\|\mathbf{w}_j + \mathbf{w}_i\|^2}{2}\right)$ .

A similar result and derivation holds for the case of  $k_{\sin, \text{Id}, \Phi}$ , which we do not reproduce here. We now mention a case that shares a connection with exponential families (see § 3.2 for a detailed description of this connection). The following and more  $\sigma = \exp$  cases are derived in Appendix C.

**Kernel 3.** Define  $\text{proj}_{\mathbb{S}^{d-1}}(\mathbf{x}) \triangleq \mathbf{x}/\|\mathbf{x}\|$  to be the projection onto the unit sphere. Then

$$k_{\text{exp}, \text{proj}_{\mathbb{S}^{d-1}}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{\Gamma(d/2) 2^{d/2-1} I_{d/2-1}(\|\mathbf{w}_i + \mathbf{w}_j\|)}{\|\mathbf{w}_i + \mathbf{w}_j\|^{d/2-1}},$$

where  $I_p$  is the modified Bessel function of the first kind of order  $p$ . In the special case  $d = 3$ , we have the closed-form  $k_{\text{exp}, \text{proj}_{\mathbb{S}^2}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{(e^{\|\mathbf{w}_i + \mathbf{w}_j\|} - e^{-\|\mathbf{w}_i + \mathbf{w}_j\|})}{2\|\mathbf{w}_i + \mathbf{w}_j\|}$ .

We end this section with a new analysis of the Snake<sub>a</sub> activation function, given by

$$\text{Snake}_a(z) = z + \frac{1}{a} \sin^2(az) = z - \frac{1}{2a} \cos(2az) + \frac{1}{2a}.$$

The Snake<sub>a</sub> function [55] is a neural network activation function that can resemble the ReLU on an interval for special choices of  $a$ , is easy to differentiate, and as we see shortly, admits certain attractive analytical tractability. We note that a similar activation function has been found using reinforcement learning to search for good activation functions [39, Table 1 and 2, row 3], up to an offset and hyperparameter  $a = 1$ . The required kernel is expressed in terms of the linear kernel (Kernel 4) and the kernel corresponding with the activation function of [39], i.e. snake without the offset,  $\text{Snake}_a(\cdot) - \frac{1}{2a}$  (Kernel 5). We first describe the linear kernel.

**Kernel 4.**  $k_{\text{Id}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \mathbf{w}_i^\top \mathbf{w}_j + b_i b_j$ .

We now derive the kernel corresponding with Snake<sub>a</sub> activation functions up to an offset.

**Kernel 5.** The kernel  $k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$  is equal to

$$\begin{aligned} & \frac{1}{4a^2} k_{\cos, \text{Id}, \Phi}(2a\boldsymbol{\theta}_i, 2a\boldsymbol{\theta}_j) + \mathbf{w}_j^\top \mathbf{w}_j \left( \sin(2ab_j) e^{-2a^2 \|\mathbf{w}_j\|^2} + \sin(2ab_i) e^{-2a^2 \|\mathbf{w}_i\|^2} \right) \\ & - \frac{b_i}{2a} \cos(2ab_j) e^{-2a^2 \|\mathbf{w}_j\|^2} - \frac{b_j}{2a} \cos(2ab_i) e^{-2a^2 \|\mathbf{w}_i\|^2} + k_{\text{Id}, \text{Id}, \Phi}(\boldsymbol{\theta}_i^{(1)}, \boldsymbol{\theta}_j^{(1)}). \end{aligned}$$

The kernel corresponding with Snake<sub>a</sub> activations is then stated in terms of Kernel 4 and 5.

**Kernel 6.** The kernel  $k_{\text{Snake}_a, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$  is equal to

$$\begin{aligned} & \frac{1}{2a} \left( b_i - \frac{1}{2a} \cos(2ab_i) \exp(-2a^2 \|\mathbf{w}_i\|^2) + b_j - \frac{1}{2a} \cos(2ab_j) \exp(-2a^2 \|\mathbf{w}_j\|^2) \right) \\ & + k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) + \frac{1}{4a^2}. \end{aligned}$$

### 3 Properties of squared neural families

#### 3.1 Fisher-Neyman factorisation and sufficient statistics

If the base measure  $\mu$  is absolutely continuous with respect to some measure  $\nu$ , and  $\frac{d\mu}{d\nu} : \Omega \rightarrow [0, \infty)$  is the Radon-Nikodym derivative, then the SNEFY admits a probability density function  $p(\cdot | \mathbf{V}, \boldsymbol{\Theta})$  with respect to  $\nu$ ,

$$p(\mathbf{x} | \mathbf{V}, \boldsymbol{\Theta}) = \underbrace{\frac{d\mu}{d\nu}(\mathbf{x})}_{\text{Independent of } \mathbf{V}, \boldsymbol{\Theta}} \times \underbrace{\frac{\|\mathbf{f}(\mathbf{t}(\mathbf{x}); \mathbf{V}, \boldsymbol{\Theta})\|_2^2}{z(\mathbf{V}, \boldsymbol{\Theta})}}_{\text{Depends on } \mathbf{x} \text{ only through } \mathbf{t}(\mathbf{x})}. \quad (7)$$

The Fisher-Neyman theorem (for example, see Theorem 6.14 of [6]) says that the existence of such a factorisation is equivalent to the fact that  $\mathbf{t}$  is a sufficient statistic for the parameters  $\mathbf{V}, \boldsymbol{\Theta}$ .

#### 3.2 Connections with exponential families

In this section we will use the activation  $\sigma(u) = \exp(u/2)$ . We note that we can absorb the bias terms  $\mathbf{b}$  into the  $\mathbf{V}$  parameters<sup>4</sup> and obtain as a special case the following family of distributions

$$P(d\mathbf{x}; \mathbf{V}, \mathbf{W}) = \frac{1}{\text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)} \sum_{i=1}^n \sum_{j=1}^n \mathbf{v}_{\cdot, i}^\top \mathbf{v}_{\cdot, j} \exp\left(\frac{1}{2}(\mathbf{w}_i + \mathbf{w}_j)^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}), \quad (8)$$

<sup>4</sup>Note that  $\exp(b_i)$  is independent of  $\mathbf{x}$  and only appears as a product with  $v_i$ .

which is a mixture<sup>5</sup> of distributions  $P_e(\cdot; \frac{1}{2}(\mathbf{w}_i + \mathbf{w}_j))$  belonging to a classical exponential family  $P_e$  [50, 32], given in the canonical form by

$$P_e(d\mathbf{x}; \mathbf{w}) = \frac{1}{z_e(\mathbf{w})} \exp(\mathbf{w}^\top \mathbf{t}(\mathbf{x})) \mu(d\mathbf{x}), \quad z_e(\mathbf{w}) = \int_{\mathbb{X}} \exp(\mathbf{w}^\top \mathbf{t}(\mathbf{x})) \mu(d\mathbf{x}). \quad (9)$$

It is helpful to identify the following three further cases:

1. When  $n = m = 1$ ,  $v_{11}^2$  cancels in the numerator and denominator and we obtain an exponential family with base measure  $\mu$  supported on  $\mathbb{X}$ , sufficient statistic  $\mathbf{t}$ , canonical parameter  $\mathbf{w}_1$  and normalising constant  $z_e(\mathbf{w}_1)$ . Every exponential family is thus a SNEFY, but not conversely.
2. When  $m > 1$  and  $n > 1$ , we obtain a type of exponential family mixture model with coefficients  $\mathbf{V}^\top \mathbf{V}$ , some of which may be negative. Advantages of allowing negative weights in mixture models in terms of learning rates are discussed in [42]. The rank of  $\mathbf{V}^\top \mathbf{V}$  is at most  $\min(m, n)$ .
3. When  $m > 1$  and  $n > 1$  and  $\mathbf{V}^\top \mathbf{V}$  is diagonal (i.e. each column in  $\mathbf{V}$  is orthogonal), there are at most  $n$  non-zero mixture coefficients, all of which are nonnegative. That is, we obtain a standard exponential family mixture model.

The kernel matrix  $\mathbf{K}_\Theta$  in the normalising constant of (8) is tractable whenever the normalising constant of the corresponding exponential family is itself tractable.

**Proposition 1.** *Denote by  $z_e(\mathbf{w})$  the normalising constant of the exponential family in (9). Then*

$$k_{\exp(\cdot/2), \mathbf{t}, \mu}(\mathbf{w}_i, \mathbf{w}_j) = z_e\left(\frac{1}{2}(\mathbf{w}_i + \mathbf{w}_j)\right). \quad (10)$$

The above kernel is well defined for any collection of  $\mathbf{w}_i$  which belong to the canonical parameter space of  $P_e$ , since the canonical parameter space is always convex [50]. This gives us a large number of tractable instances of SNEFY which correspond to exponential family mixture models allowing negative weights – a selection of examples is given in Appendix C. It is interesting that some properties of the exponential families are retained by this generalisation belonging to SNEFYs. For example, the following proposition, the proof of which is given in Appendix A, links the derivatives of the log-normalising constant to the mean and the covariance of the sufficient statistic.

**Proposition 2.** *Let  $\sigma(u) = \exp(u/2)$  and define the log-normalising constant as  $\Psi = \log z(\mathbf{V}, \Theta)$ .*

$$\text{Then } \sum_{i=1}^n \frac{\partial \Psi}{\partial \mathbf{w}_i} = \mathbb{E}[\mathbf{t}(\mathbf{x})] \quad \text{and} \quad \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \Psi}{\partial \mathbf{w}_i \mathbf{w}_j^\top} = \mathbb{E}[\mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top] - \mathbb{E}[\mathbf{t}(\mathbf{x})] \mathbb{E}[\mathbf{t}(\mathbf{x})]^\top.$$

### 3.3 Conditional distributions under SNEFY

An attractive property of SNEFY is that, under mild conditions, the family is closed under conditioning.

**Theorem 1.** *Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be jointly  $\text{SNEFY}_{\mathbb{X}_1 \times \mathbb{X}_2, \mathbf{t}, \sigma, \mu}$  with parameters  $\mathbf{V}$  and  $\Theta = ([\mathbf{W}_1, \mathbf{W}_2], \mathbf{b})$ . Assume that  $\mu(d\mathbf{x}) = \mu_1(d\mathbf{x}_1) \mu_2(d\mathbf{x}_2)$  and  $\mathbf{t}(\mathbf{x}) = (\mathbf{t}_1(\mathbf{x}_1), \mathbf{t}_2(\mathbf{x}_2))$ . Then the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2 = \mathbf{x}_2$  is  $\text{SNEFY}_{\mathbb{X}_1, \mathbf{t}_1, \sigma, \mu_1}$  with parameters  $\mathbf{V}$  and  $\Theta_{1|2} \triangleq (\mathbf{W}_1, \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b})$ .*

The proof, which we detail in Appendix A follows directly by folding the dependence on the conditioning variable  $\mathbf{x}_2$  into the bias term. We note that conditional density will typically be tractable if the joint density is tractable since they share the same activation function  $\sigma$ . Thus, whenever SNEFY corresponds to a tractable NNK with a non-zero bias, we can construct highly flexible *conditional density models* using SNEFY by taking  $\mathbf{t}_2$  itself to be a jointly trained deep neural network. Crucially,  $\mathbf{t}_2$  may be *completely unconstrained*. We use this observation in the experiments (§ 4).

<sup>5</sup>Here we are concerned with the setting where every term in the mixture model belongs to the same family, i.e. an exponential family mixture model but not a mixture of distinct exponential families.

### 3.4 Marginal distributions under SNEFY

Marginal distributions under SNEFY model for a general activation function  $\sigma$  need not belong to the same family. In the special case  $\sigma = \exp(\cdot/2)$ , SNEFY is in fact also closed under marginalisation, which we prove in Appendix D. Even in the general  $\sigma$  case, marginal distributions are tractable and admit closed forms whenever the joint SNEFY model and the conditional SNEFY are tractable.

**Theorem 2.** *Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be jointly SNEFY $_{\mathbb{X}_1 \times \mathbb{X}_2, \mathbf{t}, \sigma, \mu}$  with parameters  $\mathbf{V}$  and  $\Theta = ([\mathbf{W}_1, \mathbf{W}_2], \mathbf{b})$ . Assume that  $\mu(d\mathbf{x}) = \mu_1(d\mathbf{x}_1)\mu_2(d\mathbf{x}_2)$  and  $\mathbf{t}(\mathbf{x}) = (\mathbf{t}_1(\mathbf{x}_1), \mathbf{t}_2(\mathbf{x}_2))$ . Then the marginal distribution of  $\mathbf{x}_1$  is*

$$P_1(d\mathbf{x}_1) = \frac{\text{Tr}(\mathbf{V}^\top \mathbf{V} \tilde{\mathbf{C}}_\Theta(\mathbf{x}_1))}{z(\mathbf{V}, \Theta)} \mu_1(d\mathbf{x}_1),$$

where  $\tilde{\mathbf{C}}(\mathbf{x}_1)_{ij} = k_{\sigma, \mathbf{t}_2, \mu_2}((\mathbf{w}_{2i}, \mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i), (\mathbf{w}_{2j}, \mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j))$ .

The proof is given in Appendix A. Due to this tractability of the marginal distributions, it is straightforward to include the likelihood corresponding to incomplete observations (i.e. samples where we are missing some of the components of  $\mathbf{x}$ ) into the density estimation task.

### 3.5 Connections with kernel-based methods for nonnegative functions

SNEFY may be viewed as a neural network variant of the non-parametric kernel models for non-negative functions [25], which are constructed as follows. Let  $\psi : \mathbb{X} \rightarrow \mathbb{H}$  be a feature mapping to a (possibly infinite dimensional) Hilbert space  $\mathbb{H}$ . Let  $\mathbb{S}(\mathbb{H})$  be the set of all positive semidefinite (PSD) bounded linear operators  $\mathbf{A} : \mathbb{H} \rightarrow \mathbb{H}$ . Then

$$h_{\mathbf{A}}(\mathbf{x}) = \langle \psi(\mathbf{x}), \mathbf{A}\psi(\mathbf{x}) \rangle_{\mathbb{H}} \quad (11)$$

gives an elegant model for nonnegative functions parametrised by  $\mathbf{A} \in \mathbb{S}(\mathbb{H})$ , and their application to density modelling with respect to a base measure has also been explored [25, 42]. Note that by assuming boundedness of  $\mathbf{A}$  and  $\psi$ , the normalizing constant [25, Proposition 4] is given by  $\text{Tr}(\mathbf{A} \int_{\mathbb{X}} \psi(\mathbf{x}) \otimes \psi(\mathbf{x}) \mu(d\mathbf{x}))$ , which is analogous to our work where the normalising constant is given by  $z(\mathbf{V}, \Theta) = \text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)$ . This can be seen by replacing  $\mathbf{A}$  with  $\mathbf{V}^\top \mathbf{V}$  and replacing  $\int_{\mathbb{X}} \psi(\mathbf{x}) \otimes \psi(\mathbf{x}) \mu(d\mathbf{x})$  by  $\mathbf{K}_\Theta = \int_{\mathbb{X}} \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b}) \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b})^\top \mu(d\mathbf{x})$ .

Despite feature maps being infinite-dimensional, model (11) often reduces to an equivalent representation in finite-dimensions. [25] utilise a representer theorem when (11) is fitted to data using a regularised objective, while [42] more directly assume that the linear operator  $\mathbf{A}$  is inside the span of the features evaluated at the available data  $\{\mathbf{x}_\ell\}_{\ell=1}^N$ . The resulting model resembles SNEFY where  $n$  equals to the number  $N$  of datapoints, i.e.

$$h_{\mathbf{M}}(\mathbf{x}) = [\kappa(\mathbf{x}, \mathbf{x}_1), \dots, \kappa(\mathbf{x}, \mathbf{x}_N)]^\top \mathbf{M} [\kappa(\mathbf{x}, \mathbf{x}_1), \dots, \kappa(\mathbf{x}, \mathbf{x}_N)] \quad (12)$$

for a PSD matrix  $\mathbf{M} \in \mathbb{R}^{N \times N}$  and  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle_{\mathbb{H}}$ .

However, there are fundamental differences between (12) and SNEFY which we list below. The models can be seen as complementary and they inherit advantages and disadvantages of kernel methods and neural networks common in other settings, respectively.

- **Tractability.** Whereas we identify many tractable instances of SNEFY, the normalising constant of (12) requires computing  $\int \kappa(\mathbf{x}_i, \mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}_j) \mu(d\mathbf{x})$  – this is not generally tractable apart from some

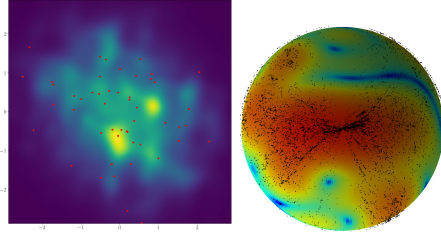


Figure 1: (Left) An instance of an (untrained) SNEFY $_{\mathbb{R}^2, \text{Id}, \text{Snake}_a, \Phi}$  density with  $n = 100$ ,  $m = 1$ ,  $v_{ij} \sim \mathcal{N}(0, 1/n)$ ,  $w_{ij} \sim \mathcal{N}(0, 4)$  and  $\mathbf{b} = \mathbf{0}$ . Shown are 50 exact samples found using rejection sampling. Numerical quadrature for this and every example supported on  $\mathbb{R}^d$  in § 4 returns a value of 1.00 for the integral over  $\mathbb{X}$ . (Right) A trained SNEFY $_{\mathbb{S}^2, \text{Id}, \text{exp}, d, \mathbf{x}}$  density with  $n = m = 30$ . Shown is the training and testing dataset [44] also used by [19] for point processes.

limited combinations of  $\kappa$  and  $\mu$  (e.g. for a Gaussian kernel  $\kappa$  and a Gaussian  $\mu$ ). Note that this kernel is evaluated at the datapoints, whereas SNEFY evaluates the kernel at the learned parameters  $w_i$ . [42] focuses on the specific case where  $\kappa$  is a Gaussian kernel, studying properties of the resulting density class which is a mixture of Gaussian densities allowing for negative weights, a model equivalent to SNEFY with the exponential activation function as described in Appendix C. Note that our treatment of SNEFY as a generalisation of the exponential family goes well beyond the Gaussian case, and that tractable instances arise with many other activation functions.

- **Expressivity.** Crucially, the feature map  $\psi$  and consequently finite dimensional representation via  $\kappa$  in (12) are treated as fixed feature maps and are not themselves learned – instead, expressivity in (12) only comes from fitting  $M$  (and potentially lengthscale hyperparameters of  $\kappa$ ), at the expense of a more involved optimisation over the space of PSD matrices. In contrast, we learn  $W$  (analogous to learning  $\psi$ ) and  $V$  jointly using neural-network style gradient optimisers. SNEFY is fully compatible with end-to-end and jointly optimised neural network frameworks, a property we leverage heavily in our experiments in § 4.
- **Conditioning.** By explicitly writing parametrisation which includes the biases, we obtain a family closed under conditioning and thus a natural model for conditional densities, whereas it is less clear how to approach conditioning when given a generic feature map  $\psi$ .

**Other related work** After submission, we became aware of another related literature, which includes mixture models with potentially negative mixture coefficients via squaring [23] and positive semi-definite probabilistic circuits [43]. We believe a marriage of ideas from SNEFY and probabilistic circuits will lead to future developments in tractable and expressive probability models.

## 4 Experiments

**Implementation** All experiments are conducted on a Dual Xeon 14-core E5-2690 with 30GB of reserved RAM and a single NVidia Tesla P100 GPU. Full experimental details are given in Appendix E. We build our implementation on top of normflows [45], a PyTorch package for normalising flows. SNEFYs are built as a BaseDistribution, which are base distributions inside a NormalizingFlow with greater than or equal to zero layers. We train all models via maximum likelihood estimation (MLE) i.e. minimising forward KL divergence.

**2D synthetic unconditional density estimation** We consider the 2 dimensional problems also benchmarked in [46]. We compare the test performance, computation time and parameter count of non-volume preserving flows (NVPs) [7] with four types of base distribution: SNEFY, resampled [46] (Resampled), diagonal Gaussian (Gauss) and Gaussian mixture model (GMM). We consider flow depths of 0, 1, 2, 4, 8 and 32 where a flow depth of 0 corresponds with the base distribution only. We use  $\sigma = \cos$ ,  $t = \text{Id}$ ,  $\mathbb{X} = \mathbb{R}^d$  and a Gaussian mixture model base density. We set  $m = 1$  and  $n = 50$ . Full architectures and further experimental details are described in Appendix E.1.

Results are shown in Table 2 for the 0 and 16 layer cases, and further tables for 1, 2, 4 and 8 layers are given in Appendix E.1. Our observations are that all base distributions are able to achieve good performance, provided they are appended with a normalising flow of appropriate depth. SNEFY is able to achieve good performance with depth 0 on all three datasets, as are Resampled and GMM for the Moons dataset. The parameter count for well-performing SNEFY models is very low, but the computation time can be relatively high. However, SNEFY is the only model which consistently achieves the highest performance within one standard deviation across all normalising flow depths.

**Data on the sphere** We compare the three cases mentioned in § 3.2 using Kernel 3, i.e. mixtures of the von Mises Fisher (VMF) distribution, as shown in Figure 1. Over 50 runs, the unconstrained  $V$ , diagonal  $V$ , and  $n = m = 1$  cases respectively obtain test negative log likelihoods of  $1.38 \pm 9.64 \times 10^{-3}$ ,  $1.46 \pm 0.016$  and  $2.34 \pm 0.26$  each in  $111.18 \pm 2.58$ ,  $109.12 \pm 0.80$  and  $61.88 \pm 0.33$  seconds (average  $\pm$  standard deviation). In this setting, allowing for a fully flexible  $V$ , going beyond the classical mixture model, shows clear benefits in performance. Results are summarised in Table 3. Full details are given in Appendix E.2.

**Conditional density estimation on astronomy data** Predicting plausible values for the velocity of distant astronomical objects (such as galaxies or quasars) without measuring their full spectra, but by



	SNEFY 0	Resampled 0	Gauss 0	GMM 0	SNEFY 16	Resampled 16	Gauss 16	GMM 16
Moons	$-1.59 \pm 0.02$	$-1.60 \pm 0.02$	$-3.29 \pm 0.02$	$-1.59 \pm 0.04$	$-1.57 \pm 0.03$	$-1.59 \pm 0.02$	$-1.68 \pm 0.24$	$-1.57 \pm 0.03$
	241	66817	4	50	19281	85857	19044	19090
Circles	$1200.84 \pm 109.77$	$541.11 \pm 12.20$	$606.35 \pm 20.29$	$629.11 \pm 14.57$	$2844.09 \pm 254.02$	$1867.26 \pm 167.75$	$2226.37 \pm 174.80$	$2276.78 \pm 211.15$
	241	66817	4	50	19281	85857	19044	19090
Rings	$643.80 \pm 150.84$	$166.35 \pm 20.59$	$52.63 \pm 3.96$	$73.52 \pm 5.47$	$2272.15 \pm 281.76$	$1533.30 \pm 158.90$	$1660.77 \pm 169.02$	$1673.54 \pm 166.84$
	241	66817	4	50	19281	85857	19044	19090
	$997.73 \pm 107.80$	$409.77 \pm 19.37$	$416.23 \pm 18.11$	$433.82 \pm 16.06$	$2654.22 \pm 269.14$	$1774.06 \pm 182.25$	$2045.65 \pm 170.80$	$2070.60 \pm 177.06$

Table 2: The first quantity is the average  $\pm$  sample standard deviation over 20 runs of test loglikelihood. The second quantity is the parameter count. The third quantity is the average  $\pm$  sample standard deviation over 20 runs of the computation time (seconds). The number in each column header is the number of non-volume preserving flow layers appended after the base distribution. Here there are 0 or 16 NVP layers. More tables showing intermediate layer results are given in Appendix E.1.

Method	Average Test NLL $\downarrow$	Compute time (s)	Parameter count
SNEFY $n = 32$	$-2.195 \pm 0.024$	$495.720 \pm 22.561$	179748
SNEFY $n = 16$	$-2.172 \pm 0.034$	$404.291 \pm 57.605$	46356
SNEFY $n = 8$	$-2.108 \pm 0.089$	$390.870 \pm 16.028$	12300
CNF $L = 4$	$-2.156 \pm 0.018$	$202.1290 \pm 10.380$	1413
CNF $L = 2$	$-2.163 \pm 0.024$	$155.090 \pm 15.809$	1155
CNF $L = 1$	$-2.171 \pm 0.012$	$122.304 \pm 1.194$	1026
CKDE	$-2.148$	$391.867$ (Not GPU-accelerated)	Train set size = $74309 \times 6$

Table 3: Performance comparison of methods on astronomy dataset. Excluding CKDE which is deterministic, an average is taken over 50 random initialisations with  $\pm$  indicating standard deviation. SNEFY shows a statistically significant average increase in performance over CNF (using a two-sample t-test), and over CKDE (using a one-sample t-test). Full experimental details in Appendix E.3.

measuring shifted spectra through broadband filters, is known as photometric redshift. Photometric redshift is an important problem in modern astronomy, as large surveys have increased the amount of available data that do not directly measure spectra. We estimate distributions for cosmological redshift  $x_1 \in \mathbb{R}$  conditional on features  $x_2 \in \mathbb{R}^5$  using a public dataset [2]. The features  $x_2$  are the magnitude a broadband filter (red) and set of four pairwise distances between broadband colour filters (ultraviolet-green, green-red, red-infrared1, infrared1-infrared2). We benchmark our own SNEFY against a conditional kernel density estimator (CKDE) and a conditional normalising flow (CNF) [54]. We use a publically available implementation [10] of CNF. Our SNEFY here is not appended with any normalising flow layers. We use  $\sigma = \text{Snake}_a$ ,  $t = \text{Id}$ ,  $\mathbb{X} = \mathbb{R}$  and a Gaussian base density. We use a deep conditional feature extractor with layer widths  $[2^5n, 2^4n, 2^3n, 2^2n, 2n]$ , and then set  $m = n/2$ . Full details are given in Appendix E.3.

While CNFs have shown great promise in modelling high dimensional image data, we expect that they are not as well-suited to nontrivial tabular data with a small to medium number of dimensions. This is because each layer is required to be an invertible transformation in the variable being modelled, so the input and output sizes must be the same, thereby significantly limiting the parameter count in each layer. We could of course increase the number of layers in the CNF model to achieve higher parameter counts, however this results in a model that is difficult to train. In our experiments, we found that increasing the depth of the CNF decreased its performance. On the other hand, SNEFY may utilise any number of parameters, as the conditioning network  $t_2$  is completely unrestricted.

**Density estimation on astronomy data using partial (marginal) observations** We perform joint probability density estimation using the redshift data. The original training dataset is a matrix in  $\mathbb{R}^{74309 \times 6}$ . We partition this matrix into batches of size 256 and randomly set each column of each batch to NaN with probability  $q$  (i.e. the corresponding dimension is missing from observations). We then train a SNEFY model that utilises partial observations by maximising the marginal likelihood according to Theorem 2. We leave the original  $74557 \times 6$  testing dataset untouched, and measure the test NLL after training. We plot the NLL as a function of  $1 - q$ , as shown in Figure 2. We also compare the performance of the SNEFY that uses partial observations with the performance of SNEFY that simply throws away incomplete observations, as well as the same normalising flow baselines that throw away incomplete observations that were used in the conditional density estimation setting. There do exist normalising flow approaches that jointly optimises conditional distributions for missing data and

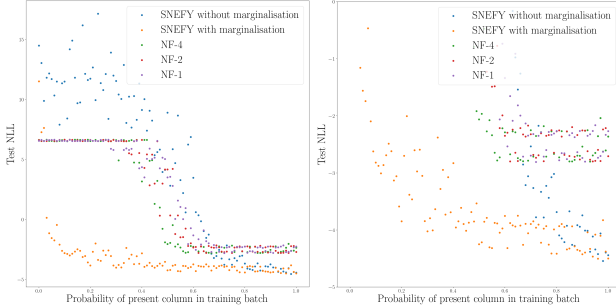


Figure 2: Density estimation under partial observations. The right plot is a zoomed in version of the left plot. NF-1, NF-2 and NF-4 are respectively normalising flows of depth 1, 2 and 4. Normalising flow models and SNEFY without marginalisation discard complete observations, whereas SNEFY use Theorem 2 to include partial observations via maximum marginal likelihood. Marginal likelihoods allow for an improved NLL.

model parameters [3] or samples from missing data and optimises for model parameters [41], however these do not allow for maximum likelihood estimation. Our observations are twofold: including partial observations improves performance, and adding more complete observations improves performance.

## 5 Discussion and conclusion

We constructed a new class of probability models – SNEFY – by normalising the squared norm of a neural network with respect to a base measure. SNEFY possesses a number of convenient properties: tractable exact normalising constants in many instances we identified, closure under conditioning, tractable marginalisation, and intriguing connections with existing models. SNEFY shows promise empirically, outperforming competing (conditional) density estimation methods in experiments.

**Sampling versus density estimation** We focus here on the problem of density estimation, for which SNEFY is well-suited. While it is sometimes possible to obtain exact SNEFY samples using rejection sampling, sampling is more computationally expensive than in other models such as normalising flows. Future work will focus on sampling, as has been done with related models [26], where the special case of Gaussian  $\sigma$  and hyperrectangular support  $\mathbb{X}$  is considered. In [26],  $r$  approximate samples with  $\mathcal{O}(r \log_2 |\mathbb{X}| + rd \log_2 \frac{2}{\rho})$  evaluations of the normalising constant are obtained. Here  $\rho$  is an approximation tolerance parameter. Note that this complexity significantly improves upon naive rejection sampling, which typically scales exponentially in dimension  $d$ . In the unconditional setting, SNEFY is constructed as only a 2-layer network, thereby limiting expressivity. However, in the conditional density estimation setting, we may use any number of layers and any architecture for the conditioning network  $t_2$ . Finally, SNEFY inherits all the usual limitations and advantages over mirroring kernel-based approaches, as discussed in § 3.5.

**Future work** We see a number of further promising future research directions. First, as we detail in Appendix F, choosing  $\sigma(\cdot) = \exp(i\cdot)$  and identity sufficient statistics results in a kernel  $k_{\exp(i\cdot), \text{Id}, \mu}$  which is the Fourier transform of a nonnegative measure. By Bochner’s theorem, the kernel is guaranteed to be (real or complex-valued) shift-invariant. The kernel matrix is Hermitian PSD (so that the normalising constant is positive and nonnegative), and we may also allow (but do not require) the readout parameters  $\mathbf{V}$  to be complex. We note that the same result would be obtained if one used a mixture of real-valued  $\cos$  and  $\sin$  activations with shared parameters (see Remark 3). Second, an alternative deep model to our deep conditional feature extractor might be to use a SNEFY model as a base measure  $\mu$  for another SNEFY model; this might be repeated  $L$  times. This leads to  $\mathcal{O}(n^{2L})$  integration terms for the normalising constant. The individual terms are tractable in certain cases, for example when  $\sigma$  is exponential or trigonometric. Third, when modelling discrete distributions with trigonometric activations, the NNK can be expressed in terms of convergent Fourier series (see Appendix G) Finally, our integration technique can be applied to other settings. For example, we may build a Poisson point process intensity function using a squared neural network and compute the intensity function in closed-form, offering a model that scales quadratically in the number of neurons  $\mathcal{O}(n^2)$  instead of comparable models which scale cubically in the number of datapoints  $\mathcal{O}(N^3)$  [9, 51].

## Acknowledgments and Disclosure of Funding

Russell and Cheng Soon would like to acknowledge the support of the Machine Learning and Artificial Intelligence Future Science Platform, CSIRO. The authors would like to thank Jia Liu for early discussions about the idea.

## References

- [1] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017.
- [2] R Beck, C-A Lin, EEO Ishida, F Gieseke, RS de Souza, MV Costa-Duarte, MW Hattab, and A Krone-Martins. On the realistic validation of photometric redshifts. *Monthly Notices of the Royal Astronomical Society*, 468(4):4323–4339, 2017.
- [3] Edgar A Bernal. Training deep normalizing flow models in highly incomplete data scenarios with prior regularization. *arXiv preprint arXiv:2104.01482*, 2021.
- [4] Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, pages 342–350, 2009.
- [5] Michael Collins, Sanjoy Dasgupta, and Robert E Schapire. A generalization of principal components analysis to the exponential family. *Advances in neural information processing systems*, 14, 2001.
- [6] Marc Peter Deisenroth, A Aldo Faisal, and Cheng Soon Ong. *Mathematics for machine learning*. Cambridge University Press, 2020.
- [7] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *International Conference on Learning Representations*, 2017.
- [8] Bradley Efron and Trevor Hastie. *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science*, volume 6. Cambridge University Press, 2021.
- [9] Seth Flaxman, Yee Whye Teh, and Dino Sejdinovic. Poisson intensity estimation with reproducing kernels. In *Artificial Intelligence and Statistics*, pages 270–279. PMLR, 2017.
- [10] Thorsten Glüsenkamp. Unifying supervised learning and vaes—automating statistical inference in (astro-) particle physics with amortized conditional normalizing flows. *arXiv e-prints*, 2020.
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- [12] Insu Han, Amir Zandieh, Jaehoon Lee, Roman Novak, Lechao Xiao, and Amin Karbasi. Fast neural kernel embeddings for general activations. In *Advances in Neural Information Processing Systems*, volume 35, pages 35657–35671, 2022.
- [13] Insu Han, Amir Zandieh, Jaehoon Lee, Roman Novak, Lechao Xiao, and Amin Karbasi. Fast neural kernel embeddings for general activations. *Advances in neural information processing systems*, 2022.
- [14] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units. *arXiv preprint arXiv:1606.08415*, 2016.
- [15] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- [16] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- [17] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.

- [18] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [19] Thomas Joseph Lawrence. Point pattern analysis on a sphere. Master’s thesis, The University of Western Australia, 2018.
- [20] Nicolas Le Roux and Yoshua Bengio. Continuous neural networks. In *Artificial Intelligence and Statistics*, pages 404–411, 2007.
- [21] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and Fugie Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- [22] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as Gaussian processes. In *International Conference on Learning Representations*, 2018.
- [23] Lorenzo Loconte, Stefan Mengel, Nicolas Gillis, and Antonio Vergari. Negative mixture models via squaring: Representation and learning. In *The 6th Workshop on Tractable Probabilistic Modeling*, 2023.
- [24] David JC MacKay. Introduction to Gaussian processes, 1998.
- [25] Ulysse Marteau-Ferey, Francis Bach, and Alessandro Rudi. Non-parametric models for non-negative functions. *Advances in neural information processing systems*, 33:12816–12826, 2020.
- [26] Ulysse Marteau-Ferey, Francis Bach, and Alessandro Rudi. Sampling from arbitrary functions via psd models. In *International Conference on Artificial Intelligence and Statistics*, pages 2823–2861. PMLR, 2022.
- [27] P. McCullagh and J.A. Nelder. *Generalized Linear Models, Second Edition*. Monographs on Statistics and Applied Probability Series. Chapman & Hall, 1989.
- [28] Lassi Meronen, Christabella Irwanto, and Arno Solin. Stationary activations for uncertainty calibration in deep learning. *Advances in Neural Information Processing Systems*, 33:2338–2350, 2020.
- [29] Lassi Meronen, Martin Trapp, and Arno Solin. Periodic activation functions induce stationarity. *Advances in Neural Information Processing Systems*, 34:1673–1685, 2021.
- [30] Charlie Nash and Conor Durkan. Autoregressive energy machines. In *International Conference on Machine Learning*, pages 1735–1744, 2019.
- [31] Radford M Neal. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995.
- [32] Frank Nielsen and Vincent Garcia. Statistical exponential families: A digest with flash cards. *arXiv preprint arXiv:0911.4863*, 2009.
- [33] Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In *The International Conference on Learning Representations*, 2019.
- [34] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research*, 22(1):2617–2680, 2021.
- [35] Georgios Papamakarios. *Neural density estimation and likelihood-free inference*. PhD thesis, University of Edinburgh, 2019.
- [36] Joouyoung Park and Irwin W Sandberg. Universal approximation using radial-basis-function networks. *Neural computation*, 3(2):246–257, 1991.

- [37] Tim Pearce, Russell Tsuchida, Mohamed Zaki, Alexandra Brintrup, and Andy Neely. Expressive priors in Bayesian neural networks: Kernel combinations and periodic functions. In *Uncertainty in Artificial Intelligence*, 2019.
- [38] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007.
- [39] Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions, 2018.
- [40] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538, 2015.
- [41] Trevor W Richardson, Wencheng Wu, Lei Lin, Beilei Xu, and Edgar A Bernal. Mcflow: Monte carlo flow models for data imputation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14205–14214, 2020.
- [42] Alessandro Rudi and Carlo Ciliberto. PSD representations for effective probability models. In *Advances in Neural Information Processing Systems*, volume 34, pages 19411–19422, 2021.
- [43] Aleksanteri Sladek. Positive Semi-Definite Probabilistic Circuits. Master’s thesis, Aalto University. School of Science, 2023.
- [44] W. Steinicke. Revised new general catalogue and index catalogue (revised ngc/ic). [http://www.klima-luft.de/steinicke/index\\_e.htm](http://www.klima-luft.de/steinicke/index_e.htm). Accessed 2nd May 2015.
- [45] Vincent Stimper, David Liu, Andrew Campbell, Vincent Berenz, Lukas Ryll, Bernhard Schölkopf, and José Miguel Hernández-Lobato. normflows: A PyTorch package for normalizing flows. *arXiv preprint arXiv:2302.12014*, 2023.
- [46] Vincent Stimper, Bernhard Schölkopf, and José Miguel Hernández-Lobato. Resampling Base Distributions of Normalizing Flows. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151, pages 4915–4936, 2022.
- [47] Russell Tsuchida, Tim Pearce, Chris van der Heide, Fred Roosta, and Marcus Gallagher. Avoiding kernel fixed points: Computing with ELU and GELU infinite networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9967–9977, 2021.
- [48] Russell Tsuchida, Fred Roosta, and Marcus Gallagher. Invariance of weight distributions in rectified MLPs. In *International Conference on Machine Learning*, pages 5002–5011, 2018.
- [49] Russell Tsuchida, Fred Roosta, and Marcus Gallagher. Richer priors for infinitely wide multi-layer perceptrons. *arXiv preprint arXiv:1911.12927*, 2019.
- [50] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [51] Christian J Walder and Adrian N Bishop. Fast bayesian intensity estimation for the permanent process. In *International Conference on Machine Learning*, pages 3579–3588. PMLR, 2017.
- [52] Christopher KI Williams. Computing with infinite networks. In *Advances in neural information processing systems*, pages 295–301, 1997.
- [53] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*. MIT press Cambridge, MA, 2006.
- [54] Christina Winkler, Daniel Worrall, Emiel Hoogeboom, and Max Welling. Learning likelihoods with conditional normalizing flows, 2019.
- [55] Liu Ziyin, Tilman Hartwig, and Masahito Ueda. Neural networks fail to learn periodic functions and how to fix it. In *Advances in Neural Information Processing Systems*, volume 33, pages 1583–1594, 2020.

## A Proofs

**Identity 1.** *The integral (2) admits a representation of the form*

$$z(\mathbf{V}, \Theta) = \text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta) \quad (6)$$

where  $k_{\sigma, \mathbf{t}, \mu}$  is as defined in (5), and  $\mathbf{K}_\Theta$  is the PSD matrix whose  $ij$ th entry is  $k_{\sigma, \mathbf{t}, \mu}(\theta_i, \theta_j)$ .

*Proof.* Let  $\widetilde{\mathbf{K}}_\Theta(\mathbf{x})$  be the PSD matrix whose  $ij$ th entry is  $\sigma(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{w}_j^\top \mathbf{t}(\mathbf{x}) + b_j)^\top$ . The squared norm of the neural network evaluation is given by  $\text{Tr}(\mathbf{V}^\top \mathbf{V} \widetilde{\mathbf{K}}_\Theta(\mathbf{x}))$ , since

$$\begin{aligned} \|\mathbf{V} \sigma(\mathbf{W} \mathbf{t}(\mathbf{x}) + \mathbf{b})\|_2^2 &= \sum_{i=1}^m \sum_{j_1=1}^n \sum_{j_2=1}^n v_{ij_1} v_{ij_2} \sigma(\mathbf{w}_{j_1}^\top \mathbf{t}(\mathbf{x}) + b_{j_1}) \sigma(\mathbf{w}_{j_2}^\top \mathbf{t}(\mathbf{x}) + b_{j_2}) \\ &= \sum_{j_1=1}^n \sum_{j_2=1}^n \mathbf{v}_{\cdot, j_1}^\top \mathbf{v}_{\cdot, j_2} \sigma(\mathbf{w}_{j_1}^\top \mathbf{t}(\mathbf{x}) + b_{j_1}) \sigma(\mathbf{w}_{j_2}^\top \mathbf{t}(\mathbf{x}) + b_{j_2}) \\ &= \langle \mathbf{V}^\top \mathbf{V}, \widetilde{\mathbf{K}}_\Theta(\mathbf{x}) \rangle_F \\ &= \text{Tr}(\mathbf{V}^\top \mathbf{V} \widetilde{\mathbf{K}}_\Theta(\mathbf{x})), \end{aligned}$$

where  $\mathbf{v}_{\cdot, j_1}$  denotes the  $j_1$ th column of  $\mathbf{V}$  and  $\langle \cdot, \cdot \rangle_F$  denotes the Frobenius inner product. Therefore using the definition (2) directly and the linearity of the Frobenius inner product, the normalising constant is

$$\begin{aligned} z(\mathbf{V}, \Theta) &= \int_{\mathbb{X}} \|\mathbf{V} \sigma(\mathbf{W} \mathbf{t}(\mathbf{x}) + \mathbf{b})\|_2^2 \mu(d\mathbf{x}) \\ &= \int_{\mathbb{X}} \text{Tr}(\mathbf{V}^\top \mathbf{V} \widetilde{\mathbf{K}}_\Theta(\mathbf{x})) \mu(d\mathbf{x}) \\ &= \text{Tr}(\mathbf{V}^\top \mathbf{V} \int_{\mathbb{X}} \widetilde{\mathbf{K}}_\Theta(\mathbf{x}) \mu(d\mathbf{x})) \\ &= \text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta). \end{aligned} \quad (13)$$

□

**Proposition 2.** *Let  $\sigma(u) = \exp(u/2)$  and define the log-normalising constant as  $\Psi = \log z(\mathbf{V}, \Theta)$ .*

$$\text{Then } \sum_{i=1}^n \frac{\partial \Psi}{\partial \mathbf{w}_i} = \mathbb{E}[\mathbf{t}(\mathbf{x})] \quad \text{and} \quad \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 \Psi}{\partial \mathbf{w}_i \partial \mathbf{w}_j^\top} = \mathbb{E}[\mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top] - \mathbb{E}[\mathbf{t}(\mathbf{x})] \mathbb{E}[\mathbf{t}(\mathbf{x})]^\top.$$

*Proof.* The result follows by noticing that the logarithmic derivative property holds,  $\sum_{i=1}^n \frac{\partial \Psi}{\partial \mathbf{w}_i} = \frac{1}{z} \sum_{i=1}^n \frac{\partial z}{\partial \mathbf{w}_i}$ , and that by writing

$$z = \sum_{i,j} \mathbf{v}_{\cdot, i}^\top \mathbf{v}_{\cdot, j} \int \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_j^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}),$$

we obtain

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{w}_i} &= \|\mathbf{v}_{\cdot, i}\|^2 \int \mathbf{t}(\mathbf{x}) \exp(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x})) \mu(d\mathbf{x}) \\ &\quad + \mathbf{v}_{\cdot, i}^\top \sum_{j \neq i} \mathbf{v}_{\cdot, j} \int \mathbf{t}(\mathbf{x}) \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_j^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}) \\ &= \mathbf{v}_{\cdot, i}^\top \sum_j \mathbf{v}_{\cdot, j} \int \mathbf{t}(\mathbf{x}) \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_j^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}). \end{aligned}$$

Now summing across  $i$  gives  $\sum_{i=1}^n \frac{\partial \Psi}{\partial \mathbf{w}_i} = \int \mathbf{t}(\mathbf{x}) P(d\mathbf{x})$  as required.

To obtain the second result, we apply the product rule to find

$$\sum_{i,j} \frac{\partial^2 \Psi}{\partial \mathbf{w}_i \mathbf{w}_j^\top} = \frac{1}{z} \sum_{i,j} \frac{\partial^2 z}{\partial \mathbf{w}_i \mathbf{w}_j^\top} - \left( \frac{1}{z} \sum_i \frac{\partial z}{\partial \mathbf{w}_i} \right) \left( \frac{1}{z} \sum_j \frac{\partial z}{\partial \mathbf{w}_j^\top} \right)$$

and note that

$$\frac{\partial^2 z}{\partial \mathbf{w}_i \mathbf{w}_j^\top} = \begin{cases} \frac{1}{2} \mathbf{v}_{\cdot,i}^\top \mathbf{v}_{\cdot,j} \int \mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_j^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}), & \text{if } i \neq j, \\ \|\mathbf{v}_{\cdot,i}\|^2 \int \mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top \exp\left(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}) \\ + \frac{1}{2} \mathbf{v}_{\cdot,i}^\top \sum_{r \neq i} \mathbf{v}_{\cdot,r} \int \mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_r^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}), & \text{if } i = j. \end{cases}$$

Thus,

$$\begin{aligned} \frac{1}{z} \sum_{i,j} \frac{\partial^2 z}{\partial \mathbf{w}_i \mathbf{w}_j^\top} &= \frac{1}{z} \sum_{i=1}^n \left( \frac{\partial^2 z}{\partial \mathbf{w}_i \mathbf{w}_i^\top} + \sum_{j \neq i} \frac{\partial^2 z}{\partial \mathbf{w}_i \mathbf{w}_j^\top} \right) \\ &= \frac{1}{z} \sum_{i,j} \mathbf{v}_i^\top \mathbf{v}_j \int \mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top \exp\left(\frac{1}{2} \mathbf{w}_i^\top \mathbf{t}(\mathbf{x})\right) \exp\left(\frac{1}{2} \mathbf{w}_j^\top \mathbf{t}(\mathbf{x})\right) \mu(d\mathbf{x}) \\ &= \int \mathbf{t}(\mathbf{x}) \mathbf{t}(\mathbf{x})^\top P(d\mathbf{x}), \end{aligned}$$

as required. □

**Remark 2.** Given the above relationship between the log-normalising constant and the expectation of the sufficient statistic, we can also ask whether the maximum likelihood estimation of the mean parameters  $\mathbb{E}[\mathbf{t}(\mathbf{x})]$  proceeds in the same way as in the exponential family case. The answer is positive but with two caveats. First, the log-likelihood need not be concave in  $\mathbf{W}$ , and may have many local optima or stationary points. Second, unlike the exponential family distribution, the SNEFY distribution is not determined by its mean parameters, so the MLE estimation of the mean parameters may not constitute a meaningful task in SNEFY modelling (unless we are in the case where precisely the expectation of  $\mathbf{t}(\mathbf{x})$  under the SNEFY model is of interest).

**Corollary 1.** Given a dataset  $\{\mathbf{x}_\ell\}_{\ell=1}^N$ , and a SNEFY model with  $\sigma(u) = \exp(\frac{1}{2}u)$ , assume that all rows of the maximum likelihood estimator of  $\mathbf{W}$  are in the interior of the natural parameter space of the corresponding exponential family. Denote the mean parameter as  $\mathbf{m} = \mathbb{E}[\mathbf{t}(\mathbf{x})]$ . Then the maximum likelihood estimate of  $\mathbf{m}$  is  $\hat{\mathbf{m}} = \frac{1}{N} \sum_{\ell=1}^N \mathbf{t}(\mathbf{x}_\ell)$ .

*Proof.* Since MLE is achieved at a stationary point of the log-likelihood, the proof follows by writing the log likelihood as

$$\sum_{\ell=1}^N \log p(\mathbf{x}_\ell; \mathbf{V}, \Theta) = \text{const} + \sum_{\ell=1}^N \log \|\mathbf{f}(\mathbf{t}(\mathbf{x}_\ell); \mathbf{V}, \Theta)\|_2^2 - N\Psi$$

and concluding that at the MLE  $\{\mathbf{w}_i^*\}_{i=1}^n$  for  $\mathbf{W}$ , we must have

$$\sum_{\ell=1}^N \frac{\partial \log \|\mathbf{f}(\mathbf{t}(\mathbf{x}_\ell); \mathbf{V}, \Theta)\|_2^2}{\partial \mathbf{w}_i^*} = N \frac{\partial \Psi}{\partial \mathbf{w}_i^*}, \quad i = 1, \dots, n. \quad (14)$$

But

$$\sum_{i=1}^n \frac{\partial \log \|\mathbf{f}(\mathbf{t}(\mathbf{x}_\ell); \mathbf{V}, \Theta)\|_2^2}{\partial \mathbf{w}_i} = \mathbf{t}(\mathbf{x}_\ell),$$

so the result follows by summing (14) over  $i$  and dividing by  $N$ . Note that maximum likelihood estimates are invariant to transformations, even if the transformation is not bijective. So if  $\{\mathbf{w}_i^*\}_{i=1}^n$  is an MLE, we may construct a mapping from  $\{\mathbf{w}_i^*\}_{i=1}^n$  to a corresponding MLE  $\hat{\mathbf{m}}$  for the mean parameter. □

**Theorem 1.** Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be jointly  $SNEFY_{\mathbb{X}_1 \times \mathbb{X}_2, \mathbf{t}, \sigma, \mu}$  with parameters  $\mathbf{V}$  and  $\Theta = ([\mathbf{W}_1, \mathbf{W}_2], \mathbf{b})$ . Assume that  $\mu(d\mathbf{x}) = \mu_1(d\mathbf{x}_1)\mu_2(d\mathbf{x}_2)$  and  $\mathbf{t}(\mathbf{x}) = (\mathbf{t}_1(\mathbf{x}_1), \mathbf{t}_2(\mathbf{x}_2))$ . Then the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2 = \mathbf{x}_2$  is  $SNEFY_{\mathbb{X}_1, \mathbf{t}_1, \sigma, \mu_1}$  with parameters  $\mathbf{V}$  and  $\Theta_{1|2} \triangleq (\mathbf{W}_1, \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b})$ .

*Proof.* The joint distribution of  $\mathbf{x}$  satisfies

$$P(d\mathbf{x}; \mathbf{V}, \Theta) \propto \left\| \mathbf{V} \sigma(\mathbf{W}_1 \mathbf{t}_1(\mathbf{x}_1) + \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b}) \right\|_2^2 \mu(d\mathbf{x}).$$

Therefore, the distribution of  $\mathbf{x}_1$  conditionally on  $\mathbf{x}_2 = \mathbf{x}_2$ , which is obtained by dividing the joint distribution by the marginal distribution of  $\mathbf{x}_2$  (which is independent of  $\mathbf{x}_1$ ), satisfies

$$P_1(d\mathbf{x}_1 | \mathbf{x}_2; \mathbf{V}, (\mathbf{W}_1, \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b})) \propto \left( \mathbf{V} \sigma(\mathbf{W}_1 \mathbf{t}_1(\mathbf{x}_1) + \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b}) \right)^2 \mu_1(d\mathbf{x}_1).$$

That is, the term  $\mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b}$  is viewed as a constant bias term when the expression on the right hand side is an unnormalised measure with respect to the variable  $\mathbf{x}_1$ .  $\square$

**Theorem 2.** Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be jointly  $SNEFY_{\mathbb{X}_1 \times \mathbb{X}_2, \mathbf{t}, \sigma, \mu}$  with parameters  $\mathbf{V}$  and  $\Theta = ([\mathbf{W}_1, \mathbf{W}_2], \mathbf{b})$ . Assume that  $\mu(d\mathbf{x}) = \mu_1(d\mathbf{x}_1)\mu_2(d\mathbf{x}_2)$  and  $\mathbf{t}(\mathbf{x}) = (\mathbf{t}_1(\mathbf{x}_1), \mathbf{t}_2(\mathbf{x}_2))$ . Then the marginal distribution of  $\mathbf{x}_1$  is

$$P_1(d\mathbf{x}_1) = \frac{\text{Tr}(\mathbf{V}^\top \mathbf{V} \tilde{\mathbf{C}}_\Theta(\mathbf{x}_1))}{z(\mathbf{V}, \Theta)} \mu_1(d\mathbf{x}_1),$$

where  $\tilde{\mathbf{C}}(\mathbf{x}_1)_{ij} = k_{\sigma, \mathbf{t}_2, \mu_2}((\mathbf{w}_{2i}, \mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i), (\mathbf{w}_{2j}, \mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j))$ .

*Proof.* The marginal distribution of the random variable  $\mathbf{x}_1$  is obtained by marginalising out the joint distribution with respect to  $\mathbf{x}_2$ ,

$$P_1(d\mathbf{x}_1) = \frac{1}{z(\mathbf{V}, \Theta)} \underbrace{\left( \int_{\mathbb{X}_2} \left\| \mathbf{V} \sigma(\mathbf{W}_1 \mathbf{t}_1(\mathbf{x}_1) + \mathbf{W}_2 \mathbf{t}_2(\mathbf{x}_2) + \mathbf{b}) \right\|_2^2 \mu_2(d\mathbf{x}_2) \right)}_{\triangleq z_2} \mu_1(d\mathbf{x}_1).$$

The integral  $z_2$  takes a similar form to  $z(\mathbf{V}, \Theta)$ ,

$$\begin{aligned} z_2 &= \text{Tr}(\mathbf{V}^\top \mathbf{V} \tilde{\mathbf{C}}_\Theta(\mathbf{x}_1)), \quad \text{where} \\ \tilde{\mathbf{C}}_{ij}(\mathbf{x}_1) &= \int_{\mathbb{X}_2} \sigma(\mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + \mathbf{w}_{2i}^\top \mathbf{t}_2(\mathbf{x}_2) + b_i) \sigma(\mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + \mathbf{w}_{2j}^\top \mathbf{t}_2(\mathbf{x}_2) + b_j) \mu_2(d\mathbf{x}_2) \\ &= k_{\sigma, \mathbf{t}_2, \mu_2}((\mathbf{w}_{2i}, \mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i), (\mathbf{w}_{2j}, \mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j)). \end{aligned}$$

$\square$

## B Derivation of neural network kernels

**Kernel 1.**  $k_{\sigma, \text{Id}, \Phi_{\mathbf{C}}, \mathbf{m}}(\theta_i, \theta_j) = k_{\sigma, \text{Id}, \Phi}(\mathcal{T}\theta_i, \mathcal{T}\theta_j)$ , where  $\mathcal{T}\Theta = (\mathbf{W}\mathbf{A}, \mathbf{b} + \mathbf{W}\mathbf{m})$ ,  $\mathcal{T}\theta_i = (\mathbf{w}_i^\top \mathbf{A}, b_i + \mathbf{w}_i^\top \mathbf{m})$  and  $\mathbf{A}$  is a matrix factor such that covariance  $\mathbf{C} = \mathbf{A}\mathbf{A}^\top$ .

*Proof.* The NNK may be expressed as an expectation with respect to a Gaussian random variable  $\mathbf{x}$  with mean  $\mathbf{m}$  and covariance matrix  $\mathbf{C}$ . It holds that  $\mathbf{x} \stackrel{d}{=} \mathbf{A}\mathbf{z} + \mathbf{m}$ , where  $\mathbf{z}$  is a zero-mean independent standard Gaussian random vector, so the kernel may be expressed in terms of an expectation over  $\mathbf{z}$  instead. More concretely,

$$\begin{aligned} k_{\sigma, \text{Id}, \Phi_{\mathbf{C}}, \mathbf{m}}(\theta_i, \theta_j) &= \mathbb{E}_{\mathbf{x}} [\sigma(\mathbf{w}_i^\top \mathbf{x} + b_i) \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)] \\ &= \mathbb{E}_{\mathbf{z}} [\sigma(\mathbf{w}_i^\top \mathbf{A}\mathbf{z} + \mathbf{w}_i^\top \mathbf{m} + b_i) \sigma(\mathbf{w}_j^\top \mathbf{A}\mathbf{z} + \mathbf{w}_j^\top \mathbf{m} + b_j)] \\ &= k_{\sigma, \text{Id}, \Phi}(\mathcal{T}\theta_i, \mathcal{T}\theta_j). \end{aligned}$$

$\square$



**Kernel 2.**  $k_{\cos, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \frac{\cos |b_i - b_j|}{2} \exp\left(-\frac{\|\mathbf{w}_j - \mathbf{w}_i\|^2}{2}\right) + \frac{\cos |b_i + b_j|}{2} \exp\left(-\frac{\|\mathbf{w}_j + \mathbf{w}_i\|^2}{2}\right).$

*Proof.* First observe that the expected value of the cosine of a Gaussian random variable can be evaluated by equating the real and imaginary components of the characteristic function of a Gaussian random variable and the expected value of Euler's form. That is, if  $z$  is Gaussian with mean  $m$  and variance  $v^2$ ,

$$\begin{aligned} \mathbb{E}e^{iz} &= \mathbb{E}[\cos(z)] + i\mathbb{E}[\sin(z)] = e^{im - \frac{1}{2}v^2} \\ &= (\cos m + i \sin m)e^{-\frac{1}{2}v^2} \\ \implies \mathbb{E}[\cos(z)] &= \cos(m)e^{-\frac{1}{2}v^2}. \end{aligned}$$

With this identity at hand, we proceed by direct evaluation of (4).

$$\begin{aligned} k_{\cos, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}_{\mathbf{x}}[\cos(\mathbf{w}_i^\top \mathbf{x} + b_i) \cos(\mathbf{w}_j^\top \mathbf{x} + b_j)], \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{x}}[\cos((\mathbf{w}_i - \mathbf{w}_j)^\top \mathbf{x} + (b_i - b_j)) + \cos((\mathbf{w}_i + \mathbf{w}_j)^\top \mathbf{x} + (b_i + b_j))] \\ &= \frac{1}{2} \cos |b_i - b_j| \exp\left(-\frac{1}{2}\|\mathbf{w}_i - \mathbf{w}_j\|^2\right) + \cos |b_i + b_j| \exp\left(-\frac{1}{2}\|\mathbf{w}_i + \mathbf{w}_j\|^2\right). \end{aligned}$$

□

**Kernel 4.**  $k_{\text{Id}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \mathbf{w}_i^\top \mathbf{w}_j + b_i b_j.$

*Proof.* This is immediate from the expected value of a product of two correlated Gaussians,  $\mathbf{w}_i^\top \mathbf{x} + b_i$  and  $\mathbf{w}_j^\top \mathbf{x} + b_j$ . □

**Kernel 5.** The kernel  $k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$  is equal to

$$\begin{aligned} &\frac{1}{4a^2} k_{\cos, \text{Id}, \Phi}(2a\boldsymbol{\theta}_i, 2a\boldsymbol{\theta}_j) + \mathbf{w}_j^\top \mathbf{w}_j \left( \sin(2ab_j) e^{-2a^2\|\mathbf{w}_j\|^2} + \sin(2ab_i) e^{-2a^2\|\mathbf{w}_i\|^2} \right) \\ &\quad - \frac{b_i}{2a} \cos(2ab_j) e^{-2a^2\|\mathbf{w}_j\|^2} - \frac{b_j}{2a} \cos(2ab_i) e^{-2a^2\|\mathbf{w}_i\|^2} + k_{\text{Id}, \text{Id}, \Phi}(\boldsymbol{\theta}_i^{(1)}, \boldsymbol{\theta}_j^{(1)}). \end{aligned}$$

*Proof.* Choosing  $\sigma = \text{Snake}_a(\cdot) - \frac{1}{2a}$  in (4), and expanding the resulting quadratic,

$$\begin{aligned} &k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) \\ &= \underbrace{\frac{1}{4a^2} \mathbb{E}_{\mathbf{x}}[\cos(2a(\mathbf{w}_i^\top \mathbf{x} + b_i)) \cos(2a(\mathbf{w}_j^\top \mathbf{x} + b_j))]}_{\text{Kernel 2}} - \underbrace{\frac{1}{2a} \mathbb{E}_{\mathbf{x}}[(\mathbf{w}_i^\top \mathbf{x} + b_i) \cos(2a(\mathbf{w}_j^\top \mathbf{x} + b_j))]}_{\triangleq A} \\ &\quad - \underbrace{\frac{1}{2a} \mathbb{E}_{\mathbf{x}}[\cos(2a(\mathbf{w}_i^\top \mathbf{x} + b_i))(\mathbf{w}_j^\top \mathbf{x} + b_j)]}_{\triangleq B} + \underbrace{\mathbb{E}_{\mathbf{x}}[(\mathbf{w}_i^\top \mathbf{x} + b_i)(\mathbf{w}_j^\top \mathbf{x} + b_j)]}_{\text{Kernel 4}}. \end{aligned} \tag{15}$$

We now evaluate  $A$  and  $B$ . The terms  $A$  and  $B$  obey a symmetry, so it suffices to evaluate term  $A$ . Term  $A$  can be evaluated using Stein's lemma,

$$\begin{aligned} A &= \frac{1}{2a} \mathbb{E}[(z_1 + b_i) \cos(z_2 + 2ab_j)], \quad (z_1, z_2)^\top \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbf{w}_i^\top \mathbf{w}_i & 2a\mathbf{w}_i^\top \mathbf{w}_j \\ 2a\mathbf{w}_j^\top \mathbf{w}_i & 4a^2\mathbf{w}_j^\top \mathbf{w}_j \end{pmatrix}\right) \\ &= \frac{1}{2a} \mathbb{E}[z_1 \cos(z_2 + 2ab_j)] + \frac{b_i}{2a} \mathbb{E}[\cos(z_2 + 2ab_j)] \\ &= -\mathbf{w}_i^\top \mathbf{w}_j \mathbb{E}[\sin(z_2 + 2ab_j)] + \frac{b_i}{2a} \mathbb{E}[\cos(z_2 + 2ab_j)] \\ &= -\mathbf{w}_i^\top \mathbf{w}_j \sin(2ab_j) \exp(-2a^2\|\mathbf{w}_j\|^2) + \frac{b_i}{2a} \cos(2ab_j) \exp(-2a^2\|\mathbf{w}_j\|^2). \end{aligned}$$

Assembling all the known individual terms in (15),

$$\begin{aligned}
& k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) \\
&= \frac{1}{4a^2} k_{\cos, \text{Id}}(2a\boldsymbol{\theta}_i, 2a\boldsymbol{\theta}_j) + \mathbf{w}_i^\top \mathbf{w}_j \left( \sin(2ab_j) \exp(-2a^2 \|\mathbf{w}_j\|^2) + \sin(2ab_i) \exp(-2a^2 \|\mathbf{w}_i\|^2) \right) \\
&\quad - \frac{b_i}{2a} \cos(2ab_j) \exp(-2a^2 \|\mathbf{w}_j\|^2) - \frac{b_j}{2a} \cos(2ab_i) \exp(-2a^2 \|\mathbf{w}_i\|^2) + k_{\text{Id}, \text{Id}}(\boldsymbol{\theta}_i^{(1)}, \boldsymbol{\theta}_j^{(1)}).
\end{aligned}$$

□

**Kernel 6.** The kernel  $k_{\text{Snake}_a, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$  is equal to

$$\begin{aligned}
& \frac{1}{2a} \left( b_i - \frac{1}{2a} \cos(2ab_i) \exp(-2a^2 \|\mathbf{w}_i\|^2) + b_j - \frac{1}{2a} \cos(2ab_j) \exp(-2a^2 \|\mathbf{w}_j\|^2) \right) \\
& + k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) + \frac{1}{4a^2}.
\end{aligned}$$

*Proof.* Choose  $\sigma = \text{Snake}_a$  and note that  $\text{Snake}_a(\cdot) = \left( \text{Snake}_a(\cdot) - \frac{1}{2a} \right) + \frac{1}{2a}$ . The Kernel 5 corresponds with the case  $\left( \text{Snake}_a(\cdot) - \frac{1}{2a} \right)$ , so we are left with three additional terms. These terms may be evaluated directly,

$$\begin{aligned}
& k_{\text{Snake}_a, \text{Id}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) \\
&= k_{\text{Snake}_a(\cdot) - \frac{1}{2a}, \text{Id}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) + \frac{1}{4a^2} + \\
& \quad \frac{1}{2a} \left( b_i - \frac{1}{2a} \cos(2ab_i) \exp(-2a^2 \|\mathbf{w}_i\|^2) + b_j - \frac{1}{2a} \cos(2ab_j) \exp(-2a^2 \|\mathbf{w}_j\|^2) \right).
\end{aligned}$$

□

## C Examples which generalise standard exponential family models

In this section, we will study examples of the SNEFY model which use activation function  $\sigma(u) = \exp(u/2)$  (or equivalently, up to scaling,  $\sigma(u) = \exp(u)$ ) and as such correspond to a notion of exponential family mixture models allowing negative weights, as discussed in Section 3.2. These examples have tractable kernels whenever the corresponding exponential family has a tractable normalising constant and we can write the kernels directly using Proposition 1.

**SNEFY Von Mises-Fisher mixtures.** The VMF distribution is a helpful way of defining the notion of a Gaussian distribution to the sphere. The following kernel may be used to define a VMF distribution. Alternatively, it may be viewed as a way of constructing a distribution supported on  $\mathbb{R}^d$  with sufficient statistics which are projected onto the sphere.

**Kernel 3.** Define  $\text{proj}_{\mathbb{S}^{d-1}}(\mathbf{x}) \triangleq \mathbf{x}/\|\mathbf{x}\|$  to be the projection onto the unit sphere. Then

$$k_{\text{exp}, \text{proj}_{\mathbb{S}^{d-1}}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{\Gamma(d/2) 2^{d/2-1} I_{d/2-1}(\|\mathbf{w}_i + \mathbf{w}_j\|)}{\|\mathbf{w}_i + \mathbf{w}_j\|^{d/2-1}},$$

where  $I_p$  is the modified Bessel function of the first kind of order  $p$ . In the special case  $d = 3$ , we have the closed-form  $k_{\text{exp}, \text{proj}_{\mathbb{S}^2}, \Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{(e^{\|\mathbf{w}_i + \mathbf{w}_j\|} - e^{-\|\mathbf{w}_i + \mathbf{w}_j\|})}{2\|\mathbf{w}_i + \mathbf{w}_j\|}$ .

*Proof.* If  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , then  $\mathbf{x}/\|\mathbf{x}\|$  is uniformly distributed on the sphere. From the normalizing constant of the von Mises-Fisher distribution, from Proposition 1, it then follows that

$$\begin{aligned}
k_{\text{exp,proj}_{\mathbb{S}^{d-1}},\Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}_{\mathbf{x}}[\exp(\mathbf{w}_i^\top \mathbf{x}/\|\mathbf{x}\| + b_i + \mathbf{w}_j^\top \mathbf{x}/\|\mathbf{x}\| + b_j)], \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
&= \exp(b_i + b_j) \int_{\mathbb{S}^{d-1}} \exp((\mathbf{w}_i + \mathbf{w}_j)^\top \mathbf{x}) d\mathbf{x} \frac{\Gamma(d/2)}{2\pi^{d/2}} \\
&= \frac{\exp(b_i + b_j)\Gamma(d/2)}{2\pi^{d/2}} \int_{\mathbb{S}^{d-1}} \exp(\|\mathbf{w}_i + \mathbf{w}_j\| \mathbf{a}^\top \mathbf{x}) d\mathbf{x}, \quad \text{where } \mathbf{a} \text{ is a unit vector} \\
&= \frac{\exp(b_i + b_j)\Gamma(d/2)}{2\pi^{d/2}} \frac{(2\pi)^{d/2} I_{d/2-1}(\|\mathbf{w}_i + \mathbf{w}_j\|)}{\|\mathbf{w}_i + \mathbf{w}_j\|^{d/2-1}} \\
&= \exp(b_i + b_j) \frac{\Gamma(d/2) 2^{d/2-1} I_{d/2-1}(\|\mathbf{w}_i + \mathbf{w}_j\|)}{\|\mathbf{w}_i + \mathbf{w}_j\|^{d/2-1}},
\end{aligned}$$

where  $I_p$  is the modified Bessel function of the first kind of order  $p$ . In the special case of  $p = 1/2$ , we have  $I_{1/2}(z) = \sqrt{\frac{2}{\pi z}} \sinh(z) = (\exp(z) - \exp(-z))\sqrt{\frac{1}{2\pi z}}$ . This implies that when  $d = 3$ , since  $\Gamma(3/2) = \frac{\sqrt{\pi}}{2}$ ,

$$k_{\text{exp,proj}_{\mathbb{S}^2},\Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{(e^{\|\mathbf{w}_i + \mathbf{w}_j\|} - e^{-\|\mathbf{w}_i + \mathbf{w}_j\|})}{2\|\mathbf{w}_i + \mathbf{w}_j\|}.$$

□

Note that  $k_{\text{exp,proj}_{\mathbb{S}^{d-1}},\Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = k_{\text{exp,Id},\nu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$ , where  $\nu$  is the uniform measure on the sphere  $\mathbb{S}^{d-1}$ , because if  $\mathbf{x}$  is Gaussian then  $\mathbf{x}/\|\mathbf{x}\|$  is uniform on the sphere. This allows one to construct SNEFY<sub>exp,Id, $\nu$</sub>  distributions, which are certain ‘‘mixtures’’ of VMF distributions with weights  $\mathbf{V}^\top \mathbf{V}$ .

**SNEFY Gaussian mixtures, fixed variance.** We may similarly construct kernels corresponding to ‘‘mixtures’’ of Gaussian distributions. The case here corresponds to a case of known fixed variance parameter.

**Kernel 7.**  $k_{\text{exp,Id},\Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \exp(\frac{1}{2}\|\mathbf{w}_i + \mathbf{w}_j\|^2)$ .

*Proof.* This is a consequence of the moment generating function of the multivariate Gaussian distribution. More concretely, by Proposition 1,

$$\begin{aligned}
k_{\text{exp,Id},\Phi}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}_{\mathbf{x}}[\exp(\mathbf{w}_i^\top \mathbf{x} + b_i + \mathbf{w}_j^\top \mathbf{x} + b_j)], \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\
&= \exp(b_i + b_j) \mathbb{E}_{\mathbf{x}}[\exp((\mathbf{w}_i + \mathbf{w}_j)^\top \mathbf{x})] \\
&= \exp(b_i + b_j) \exp(\frac{1}{2}\|\mathbf{w}_i + \mathbf{w}_j\|^2).
\end{aligned}$$

□

**SNEFY Poisson mixtures.** Most of our examples deal with continuous distributions but in fact SNEFY can readily be used for discrete distribution modelling. This is particularly helpful when the support is large or infinite, for which computing normalising constants can naively be challenging even in the discrete setting. Let  $\mathbb{X} = \{0, 1, 2, \dots\}$ ,  $t(x) = x$ , and the base measure  $\mu(dx) = 1/x! \nu(dx)$ , where  $\nu$  is the counting measure. A SNEFY model for a probability mass function which is a mixture of Poisson distributions allowing negative weights is given by

$$\begin{aligned}
p(x; \mathbf{V}, \mathbf{w}) &= \frac{1}{\text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)} \frac{1}{x!} \sum_{i=1}^n \sum_{j=1}^n \mathbf{v}_{:,i}^\top \mathbf{v}_{:,j} \exp((w_i + w_j)x) \\
&= \frac{1}{\text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)} \frac{1}{x!} \sum_{i=1}^n \sum_{j=1}^n \mathbf{v}_{:,i}^\top \mathbf{v}_{:,j} (\lambda_i \lambda_j)^x, \quad x = 0, 1, 2, \dots
\end{aligned}$$

following the usual mean parametrisation  $\lambda_i = e^{w_i}$ , so the individual mixture components have rates which are geometric means of  $(\lambda_i, \lambda_j)$  pairs.

**Kernel 8.** Choose the base measure  $\mu(dx) = (x!)^{-1} \nu(dx)$ , where  $\nu$  is the counting measure. We have

$$k_{\text{exp,Id},(x!)^{-1}\nu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \exp(\exp(w_i + w_j)).$$

*Proof.* This is again direct from Proposition 1. In detail,

$$\begin{aligned} k_{\text{exp,Id},(x!)^{-1}\nu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \sum_{x=0}^{\infty} \frac{1}{x!} \exp(w_i x + w_j x + b_i + b_j) \\ &= \exp(b_i + b_j) \sum_{x=0}^{\infty} \frac{1}{x!} \exp(w_i x + w_j x) \end{aligned}$$

The second factor involving the sum is the partition function of the Poisson distribution in canonical form, which is  $\exp(\exp(w_i + w_j))$ .  $\square$

The usual mean parameterisation of the Poisson distribution is through a rate parameter  $\lambda_i = \exp(w_i)$ , which would lead to the kernel representation

$$k_{\text{exp,Id},(x!)^{-1}\nu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \exp(\lambda_i \lambda_j).$$

**SNEFY Gaussian mixtures, unknown variance (Squared radial basis function network).** We now discuss an intriguing connection between the Gaussian distribution and squared RBF networks. This connection is made possible through our machinery of SNEFY distributions. Let  $\mathbb{X} = \mathbb{R}^d$ ,  $\mathbf{t}(\mathbf{x}) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2)$  (i.e.  $D = 2d$ ) and suppose  $\mu(d\mathbf{x}) = d\mathbf{x}$  is Lebesgue measure. Choose  $\sigma(\cdot) = \exp(\cdot/2)$  and consider the  $r$ -th output of our network  $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^m$

$$f_r(\mathbf{t}(\mathbf{x}); \mathbf{V}, \boldsymbol{\Theta}) = \sum_{i=1}^n v_{ri} \exp\left(\frac{1}{2} \left( \sum_{\ell=1}^d w_{i\ell} x_\ell + \sum_{\ell=1}^d \tilde{w}_{i\ell} x_\ell^2 \right)\right),$$

where we denoted  $\tilde{w}_{i\ell} = w_{i,d+\ell}$ . In this case, we require that  $\tilde{w}_{i\ell} < 0$  for the model to be (square) integrable. Reparametrising  $\sigma_{i\ell}^2 = -\frac{1}{2\tilde{w}_{i\ell}}$  and  $\mu_{i\ell} = -\frac{w_{i\ell}}{2\tilde{w}_{i\ell}}$  and absorbing the factor  $\exp\left(-\sum_{\ell=1}^d \frac{\mu_{i\ell}^2}{4\sigma_{i\ell}^2}\right)$  into readout parameters  $\mathbf{V}$ , gives

$$f_r(\mathbf{t}(\mathbf{x}); \mathbf{V}, \boldsymbol{\Theta}) = \sum_{i=1}^n v_{ri} \exp\left(-\sum_{\ell=1}^d \frac{(x_\ell - \mu_{i\ell})^2}{4\sigma_{i\ell}^2}\right).$$

Thus, we have recovered a classical radial basis function (RBF) network [52]. These models are well known to have universal approximation properties [36]. For the most commonly used form of the RBF network, we can restrict the parameters  $\sigma_{i\ell}^2 = \sigma_i^2$  to be the same across the dimensions, giving

$$f_r(\mathbf{t}(\mathbf{x}); \mathbf{V}, \boldsymbol{\Theta}) = \sum_{i=1}^n v_{ri} \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{4\sigma_i^2}\right), \quad \mathbf{x} \in \mathbb{R}^d,$$

with location parameters  $\boldsymbol{\mu}_i \in \mathbb{R}^d$  and the scale parameters  $\sigma_i^2 > 0$ . Note the unusual factor of 4 in front of  $\sigma_i^2$  – this ensures that our model in fact reduces to the usual parametrisation of multivariate normal densities, since we will be modelling the density using the squared norm of  $\mathbf{f}$ .

Since  $\mu$  is the Lebesgue measure, SNEFY gives us a density model with respect to the Lebesgue measure as

$$p(\mathbf{x}; \mathbf{V}, \boldsymbol{\Theta}) = \frac{1}{\text{Tr}(\mathbf{V}^\top \mathbf{V} \mathbf{K}_\Theta)} \sum_{i=1}^n \sum_{j=1}^n \mathbf{v}_{\cdot,i}^\top \mathbf{v}_{\cdot,j} \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{4\sigma_i^2}\right) \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_j\|^2}{4\sigma_j^2}\right).$$

If  $n = 1$ , we recover simply a multivariate normal density  $\mathcal{N}(\boldsymbol{\mu}_1, \sigma_1^2 I)$ . The above model is essentially the same as the one in [42], despite being derived in a very different way.

**Kernel 9.** Let  $\mathbf{t}^{(2)}(\mathbf{x}) = (x_1, \dots, x_d, x_1^2, \dots, x_d^2)$  so that  $D = 2d$  be the sufficient statistic. Partition  $\mathbf{W} = [\mathbf{W}_{[:,1:d]}, \tilde{\mathbf{W}}]$  and suppose  $\tilde{\mathbf{W}} < 0$  element-wise. Choose  $\mu(d\mathbf{x}) = d\mathbf{x}$  to be the Lebesgue measure. Then

$$k_{\exp, \mathbf{t}^{(2)}, d\mathbf{x}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \pi^{d/2} \exp(b_i + b_j) \prod_{l=1}^d \exp\left(-\frac{(w_{il} + w_{jl})^2}{4(\tilde{w}_{il} + \tilde{w}_{jl})}\right) \frac{1}{\sqrt{-(\tilde{w}_{il} + \tilde{w}_{jl})}}$$

*Proof.* As with the kernels above, this follows from Proposition 1, since

$$\begin{aligned} k_{\exp, \mathbf{t}^{(2)}, d\mathbf{x}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \int_{\mathbb{R}^d} \exp\left(\mathbf{w}_{i,1:d}^\top \mathbf{x} + \tilde{\mathbf{w}}_i^\top \mathbf{x}^2 + b_i + \mathbf{w}_{j,1:d}^\top \mathbf{x} + \tilde{\mathbf{w}}_j^\top \mathbf{x}^2 + b_j\right) d\mathbf{x} \\ &= (2\pi)^{d/2} \exp(b_i + b_j) \prod_{l=1}^d \exp\left(-\frac{(w_{il} + w_{jl})^2}{4(\tilde{w}_{il} + \tilde{w}_{jl})}\right) \frac{1}{\sqrt{-2(\tilde{w}_{il} + \tilde{w}_{jl})}}. \end{aligned}$$

□

While Proposition 1 gives us an expression for the kernel matrix  $\mathbf{K}_\Theta$  in terms of natural parameters, we can also express it directly in terms of parameters  $\boldsymbol{\mu}_i, \sigma_i^2$ . In particular,

$$\begin{aligned} [\mathbf{K}_\Theta]_{ij} &= \exp(b_i + b_j) \int_{\mathbb{R}^d} \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{4\sigma_i^2}\right) \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_j\|^2}{4\sigma_j^2}\right) d\mathbf{x} \\ &= \exp(b_i + b_j) \left(\frac{4\pi\sigma_i^2\sigma_j^2}{\sigma_i^2 + \sigma_j^2}\right)^{d/2} \exp\left(-\frac{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2}{4(\sigma_i^2 + \sigma_j^2)}\right). \end{aligned}$$

We briefly state two more cases without an extended discussion.

**SNEFY Gamma mixtures.**

**Kernel 10.** Let  $\mathbb{X} = (0, \infty)$ ,  $\mathbf{t}(x) = (\log x, -x)$  and  $\sigma = \exp$ . Partition  $\mathbf{W} = [\mathbf{W}_{[:,1:d]}, \tilde{\mathbf{W}}]$  and suppose  $\mathbf{W}_{[:,1:d]} > -1$  and  $\tilde{\mathbf{W}} > 0$  element-wise. Choose  $\mu(d\mathbf{x}) = d\mathbf{x}$  to be the Lebesgue measure. Then

$$k_{\exp, \mathbf{t}, d\mathbf{x}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{\Gamma(w_{i1} + w_{j1} + 1)}{(w_{i2} + w_{j2})^{w_{i1} + w_{j1} + 1}}.$$

**SNEFY Dirichlet mixtures.**

**Kernel 11.** Let  $\mathbb{X} = \Delta^{D-1}$ , a  $(D-1)$ -simplex of probability distributions, i.e.  $\mathbf{x} \in [0, 1]^D$ ,  $\sum_{i=1}^D x_i = 1$ . Let  $\sigma = \exp$ . Let  $\mathbf{t}(\mathbf{x}) = (\log x_1, \dots, \log x_D)$ . Choose  $\mu(d\mathbf{x}) = d\mathbf{x}$  to be the Lebesgue measure. Suppose  $\tilde{\mathbf{W}} > -1$  elementwise. Then

$$k_{\exp, \mathbf{t}, d\mathbf{x}}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp(b_i + b_j) \frac{\prod_{d=1}^D \Gamma(w_{id} + w_{jd} + 1)}{\Gamma\left(D + \sum_{d=1}^D w_{id} + w_{jd}\right)}.$$

## D Marginalisation in the case $\sigma = \exp(\cdot/2)$

Let  $\mathbf{M}$  be a positive semi-definite  $n \times n$  matrix. We will make use of the following SNEFY parametrisation

$$P(d\mathbf{x}; \mathbf{M}, \Theta) = \frac{1}{z(\mathbf{M}, \Theta)} \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b})^\top \mathbf{M} \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b}) \mu(d\mathbf{x}). \quad (16)$$

Since we can always write  $\mathbf{M} = \mathbf{V}^\top \mathbf{V}$ , for an  $m \times n$  matrix  $\mathbf{V}$ ,  $m \leq n$ , we have

$$\sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b})^\top \mathbf{M} \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b}) = \|\mathbf{V} \sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b})\|^2 = \sum_{i=1}^m \left( \sum_{j=1}^n v_{ij} \sigma(w_j^\top \mathbf{t}(\mathbf{x}) + b_j) \right)^2,$$

which is, as in the parametrisation given in the main text, simply the squared Euclidean norm of a multi-output neural network, with  $\mathbf{V}$  corresponding to the weights of the second layer. If we denote by  $\mathbf{v}_{\cdot j}$  the  $j$ -th column of  $\mathbf{V}$ , the normalizing constant is given by

$$\sum_{j=1}^n \sum_{l=1}^n \mathbf{v}_{\cdot j}^\top \mathbf{v}_{\cdot l} k_{\sigma, \mathbf{t}, \mu}(\boldsymbol{\theta}_j, \boldsymbol{\theta}_l) = \sum_{j=1}^n \sum_{l=1}^n m_{jl} k_{\sigma, \mathbf{t}, \mu}(\boldsymbol{\theta}_j, \boldsymbol{\theta}_l) = \text{Tr}(\mathbf{M}\mathbf{K}_\Theta)$$

where as before

$$k_{\sigma, \mathbf{t}, \mu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \int \sigma(\mathbf{w}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{w}_j^\top \mathbf{t}(\mathbf{x}) + b_j) \mu(d\mathbf{x}).$$

Now if we let  $\sigma = \exp(\cdot/2)$  we obtain a family which is also closed under marginalisation (in addition to conditioning). The Proposition below generalises Proposition 1 of [42], which considers the special case of the Gaussian PSD mixtures.

**Proposition 3.** *Let  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  be jointly  $\text{SNEFY}_{\mathbb{X}_1 \times \mathbb{X}_2, \mathbf{t}, \exp(\cdot/2), \mu}$  with parameters  $\mathbf{V}$  and  $\Theta = ([\mathbf{W}_1, \mathbf{W}_2], \mathbf{b})$ . Assume that  $\mu(d\mathbf{x}) = \mu_1(d\mathbf{x}_1)\mu_2(d\mathbf{x}_2)$  and  $\mathbf{t}(\mathbf{x}) = (\mathbf{t}_1(\mathbf{x}_1), \mathbf{t}_2(\mathbf{x}_2))$ . Then the marginal distribution of  $\mathbf{x}_1$  is  $\text{SNEFY}_{\mathbb{X}_1, \mathbf{t}_1, \exp(\cdot/2), \mu_1}$  with parameters  $\tilde{\mathbf{V}}$  and  $\Theta = (\mathbf{W}_1, \mathbf{b})$ , for some matrix  $\tilde{\mathbf{V}} \in \mathbb{R}^{m \times n}$ .*

*Proof.* Proposition 1 gives us the normalising constant for the parametrisation where biases are absorbed into  $\mathbf{V}$ . If we explicitly keep the biases in the parametrisation, we have

$$k_{\exp(\cdot/2), \mathbf{t}, \mu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \exp\left(\frac{1}{2}(b_i + b_j)\right) z_e\left(\frac{1}{2}(\mathbf{w}_i + \mathbf{w}_j)\right), \quad (17)$$

where  $z_e$  is the normalizing constant of the exponential family with the sufficient statistic  $\mathbf{t}$  and base measure  $\mu$ . By Theorem 2, we have that

$$P_1(d\mathbf{x}_1; \mathbf{M}, \Theta) = \frac{\text{Tr}(\mathbf{M}\mathbf{C}_\Theta(\mathbf{x}_1))}{\text{Tr}(\mathbf{M}\mathbf{K}_\Theta)} \mu_1(d\mathbf{x}_1),$$

where  $[\mathbf{C}_\Theta(\mathbf{x}_1)]_{ij} = k_{\sigma, \mathbf{t}_2, \mu_2}\left((\mathbf{w}_{2i}, \mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i), (\mathbf{w}_{2j}, \mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j)\right)$ . But now since  $\sigma = \exp(\cdot/2)$ , applying (17) gives

$$[\mathbf{C}_\Theta(\mathbf{x}_1)]_{ij} = z_{e,2}\left(\frac{1}{2}(\mathbf{w}_{2i} + \mathbf{w}_{2j})\right) \exp\left(\frac{1}{2}(\mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i)\right) \exp\left(\frac{1}{2}(\mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j)\right),$$

where  $z_{e,2}$  is the normalizing constant of the exponential family with the sufficient statistic  $\mathbf{t}_2$  and base measure  $\mu_2$ . Thus, we can write

$$\begin{aligned} \text{Tr}(\mathbf{M}\mathbf{C}_\Theta(\mathbf{x}_1)) &= \sum_{i=1}^n \sum_{j=1}^n \left\{ m_{ij} z_{e,2}\left(\frac{1}{2}(\mathbf{w}_{2i} + \mathbf{w}_{2j})\right) \right. \\ &\quad \left. \cdot \exp\left(\frac{1}{2}(\mathbf{w}_{1i}^\top \mathbf{t}_1(\mathbf{x}_1) + b_i)\right) \exp\left(\frac{1}{2}(\mathbf{w}_{1j}^\top \mathbf{t}_1(\mathbf{x}_1) + b_j)\right) \right\} \\ &= \exp\left(\frac{1}{2}(\mathbf{W}_1^\top \mathbf{t}_1(\mathbf{x}_1) + \mathbf{b})\right)^\top \tilde{\mathbf{M}} \exp\left(\frac{1}{2}(\mathbf{W}_1^\top \mathbf{t}_1(\mathbf{x}_1) + \mathbf{b})\right) \end{aligned}$$

and we conclude that the marginal is in the same family with  $\tilde{\mathbf{M}} = \mathbf{M} \circ \mathbf{Z}_{e,2}$ , where

$$[\mathbf{Z}_{e,2}]_{i,j} = z_{e,2}\left(\frac{1}{2}(\mathbf{w}_{2i} + \mathbf{w}_{2j})\right).$$

Note that  $\tilde{\mathbf{M}}$  is PSD as an Hadamard product of two PSD matrices. Thus, we can find  $\tilde{\mathbf{V}}$  such that  $\tilde{\mathbf{M}} = \tilde{\mathbf{V}}^\top \tilde{\mathbf{V}}$ .  $\square$

	SNEFY 1	Resampled 1	Gauss 1	GMM 1
Moons	$-1.59 \pm 0.03$ 1431 $1352.19 \pm 134.30$	$-1.76 \pm 0.02$ 68007 $615.57 \pm 24.35$	$-3.29 \pm 0.01$ 1194 $730.21 \pm 23.60$	$-1.57 \pm 0.03$ 1240 $753.16 \pm 31.41$
Circles	$-1.92 \pm 0.02$ 1431 $798.13 \pm 196.89$	$-1.94 \pm 0.03$ 68007 $291.03 \pm 29.46$	$-3.37 \pm 0.01$ 1194 $162.64 \pm 17.19$	$-2.00 \pm 0.06$ 1240 $186.65 \pm 21.10$
Rings	$-2.35 \pm 0.04$ 1431 $1134.90 \pm 133.10$	$-2.31 \pm 0.02$ 68007 $502.32 \pm 24.51$	$-3.16 \pm 0.01$ 1194 $528.50 \pm 24.51$	$-2.43 \pm 0.04$ 1240 $559.12 \pm 30.90$

Table 4: As in Table 2, but with 1 NVP layer.

	SNEFY 2	Resampled 2	Gauss 2	GMM 2
Moons	$-1.58 \pm 0.02$ 2621 $1453.05 \pm 147.67$	$-1.59 \pm 0.03$ 69197 $700.98 \pm 32.94$	$-1.62 \pm 0.03$ 2384 $831.40 \pm 49.24$	$-1.58 \pm 0.02$ 2430 $855.04 \pm 53.01$
Circles	$-1.91 \pm 0.04$ 2621 $899.77 \pm 196.54$	$-2.07 \pm 0.04$ 69197 $363.91 \pm 36.62$	$-2.66 \pm 0.05$ 2384 $267.23 \pm 32.05$	$-1.96 \pm 0.04$ 2430 $289.96 \pm 33.71$
Rings	$-2.35 \pm 0.04$ 2621 $1255.92 \pm 168.80$	$-2.32 \pm 0.03$ 69197 $579.02 \pm 23.70$	$-2.80 \pm 0.02$ 2384 $637.73 \pm 34.04$	$-2.36 \pm 0.03$ 2430 $658.98 \pm 37.69$

Table 5: As in Table 2, but with 2 NVP layers.

## E Experiments

### E.1 2D Unconditional density estimation

Our benchmarking protocol is slightly altered compared with [46]. Firstly, we measure performance over 20 random seeds instead of 1 fixed seed. We find that sometimes the variance over random seeds can be large (e.g. Resampled 0 on Circles). Secondly, rather than computing test performance at the last epoch of training, we follow the more standard procedure of returning the test performance evaluated at the epoch corresponding with the smallest validation performance. This validation/test monitoring results in substantial performance gains in all of the models, with no extra computational cost for SNEFY, Gauss and GMM. For Resampled, monitoring the validation performance to a high precision requires estimating the normalising constant to a high precision, which is computationally challenging. We therefore only check the validation performance every 100 epochs, and compute a high precision normalising constant if the validation performance is the lowest encountered so far. We train each models for a maximum of 20000 iterations (while monitoring validation performance). We use Adam with default hyperparameters and weight decay  $10^{-3}$ . The batch size is  $2^{10}$ .

We use a SNEFY with Gaussian mixture model base measure supported on  $\mathbb{X} = \mathbb{R}^2$ , identity sufficient statistic  $t$  and activation function  $\cos$ . The base measure consists of 8 mixture components, each with a diagonal covariance matrix. We use the same Resampled architecture as in the original paper [46]. We use an MLP with layer widths  $[2, 256, 256, 1]$  and sigmoid activations for the resampling distribution. We use a discount factor of 0.1 for the exponential moving average partition function calculation. Note that for Resampled models we are only able to provide an estimate of the test log likelihood, not the exact log likelihood as in other methods. We choose  $n = 50$  and  $m = 1$ . The Gaussian mixture model has 10 mixture components, each with a diagonal covariance matrix. Each normalising flow block consists of an affine coupling block, a permutation layer, and an actnorm layer.

### E.2 Modelling distributions on the sphere

We compare the performance of a VMF distribution, a “regular” VMF mixture distribution and our SNEFY construction of the VMF mixture, which allows for some negative weight coefficients as discussed in § 3.2. We use the dataset [44] also used by [19] in a different context for point

	SNEFY 4	Resampled 4	Gauss 4	GMM 4
Moons	$-1.58 \pm 0.03$	$-1.60 \pm 0.07$	$-1.60 \pm 0.03$	$-1.57 \pm 0.03$
	5001	71577	4764	4810
	$1616.21 \pm 152.61$	$851.29 \pm 33.63$	$1005.16 \pm 30.59$	$1026.76 \pm 28.21$
Circles	$-1.92 \pm 0.04$	$-1.99 \pm 0.03$	$-2.14 \pm 0.16$	$-1.94 \pm 0.04$
	5001	71577	4764	4810
	$1035.64 \pm 142.39$	$514.02 \pm 47.98$	$455.58 \pm 41.19$	$480.15 \pm 44.61$
Rings	$-2.33 \pm 0.02$	$-2.31 \pm 0.03$	$-2.59 \pm 0.12$	$-2.32 \pm 0.04$
	5001	71577	4764	4810
	$1410.36 \pm 141.04$	$749.10 \pm 52.05$	$813.92 \pm 27.20$	$833.36 \pm 30.77$

Table 6: As in Table 2, but with 4 NVP layers.

	SNEFY 8	Resampled 8	Gauss 8	GMM 8
Moons	$-1.60 \pm 0.02$	$-1.58 \pm 0.02$	$-1.61 \pm 0.08$	$-1.58 \pm 0.03$
	9761	76337	9524	9570
	$2036.15 \pm 205.52$	$1198.04 \pm 102.71$	$1399.74 \pm 65.78$	$1423.17 \pm 82.84$
Circles	$-1.91 \pm 0.03$	$-1.97 \pm 0.05$	$-2.16 \pm 0.29$	$-1.93 \pm 0.03$
	9761	76337	9524	9570
	$1405.27 \pm 144.72$	$839.23 \pm 92.99$	$838.62 \pm 82.98$	$862.25 \pm 81.09$
Rings	$-2.34 \pm 0.06$	$-2.37 \pm 0.17$	$-2.49 \pm 0.16$	$-2.33 \pm 0.07$
	9761	76337	9524	9570
	$1795.38 \pm 186.62$	$1071.50 \pm 109.12$	$1201.29 \pm 71.76$	$1218.31 \pm 70.36$

Table 7: As in Table 2, but with 8 NVP layers.

processes. This dataset, retrieved in 2015, is called the ‘‘Revised New General Catalogue and Index Catalogue’’ (RNGC/IC). The RNGC/IC consists of locations of some 10610 galaxies. We use the spherical coordinates of these galaxies and map them to the surface of a sphere.

We compare two types of  $\text{SNEFY}_{\mathbb{S}^2, \text{exp}, \text{Id}, d\mathbf{x}}$  with  $m = n = 30$ : one is constrained so that  $\mathbf{V}$  is diagonal (and therefore  $\mathbf{V}^\top \mathbf{V}$  is diagonal with nonnegative entries), and the other uses a unconstrained general  $\mathbf{V}$ . We expect the unconstrained model to be more expressive, and therefore obtain a better test NLL.

We use Adam with default hyperparameters and a full batch size. We randomly shuffle the data and perform an 80/20 train/test split. Each of the 50 runs uses a randomly sampled initialisation and train/test split. We train for 20000 epochs.

### E.3 Conditional density estimation of photometric redshift

Our deep conditional feature extractor is an MLP that uses ReLU activations in each layer and batch normalisation between each layer. Our SNEFY model uses Snake<sub>a</sub> activations with  $a = 10$ . The CNF models utilise affine layers with tanh activations. For SNEFY and CNF, we preprocess the input features so that they have sample mean zero and sample variance one. We train both deep learning models for 100 epochs using Adam with default hyperparameters and a batch size of 256. The dataset consists of 74309 training and 74557 examples, which are fixed between each run. Each run uses a randomly sampled initialisation (except for CKDE, which is deterministic).

## F Complex activation case

Here we consider an extension where we allow the neural network activation  $\sigma$  to be complex-valued, and accordingly the readout parameters  $\mathbf{V}$  to be complex, i.e.  $\mathbf{V} \in \mathbb{C}^{m \times n}$ . Note that in this case

$$\|\mathbf{f}(\mathbf{t}(\mathbf{x}); \mathbf{V}, \Theta)\|_2^2 = \|\mathbf{V}\sigma(\mathbf{W}\mathbf{t}(\mathbf{x}) + \mathbf{b})\|^2 = \sum_{r=1}^m \left| \sum_{j=1}^n v_{rj} \sigma(\mathbf{w}_j^\top \mathbf{t}(\mathbf{x}) + b_j) \right|^2.$$



Take  $\sigma(u) = \exp(iu)$  and  $\mathbf{t}(\mathbf{x}) = \mathbf{x}$ .

Then the above takes the form

$$\begin{aligned} \sum_{r=1}^m \left| \sum_{j=1}^n v_{rj} \exp(i(\mathbf{w}_j^\top \mathbf{x} + b_j)) \right|^2 &= \sum_{r=1}^m \sum_{j=1}^n \sum_{l=1}^n v_{rj} \bar{v}_{rl} e^{i(\mathbf{w}_j^\top \mathbf{x} + b_j)} e^{-i(\mathbf{w}_l^\top \mathbf{x} + b_l)} \\ &= \sum_{r=1}^m \sum_{j=1}^n \sum_{l=1}^n v_{rj} \bar{v}_{rl} e^{i(b_j - b_l)} e^{i((\mathbf{w}_j - \mathbf{w}_l)^\top \mathbf{x})}. \end{aligned}$$

As in the  $\sigma = \exp$  case, bias terms can be folded into the readout parameters  $\mathbf{V}$  (which is the reason why we also require readout parameters to be complex). We note that the case  $n = 1$  is not interesting as it simply reduces to the (normalised) base measure  $\mu$ .

In order to obtain the normalizing constant  $\text{Tr}(\mathbf{V}^H \mathbf{V} \mathbf{K}_\Theta)$ , we need the integral of the form

$$[\mathbf{K}_\Theta]_{jl} = k_{\exp(i\cdot), \text{Id}, \mu}(\mathbf{w}_j, \mathbf{w}_l) = \int e^{i((\mathbf{w}_j - \mathbf{w}_l)^\top \mathbf{x})} \mu(d\mathbf{x}) =: \kappa(\mathbf{w}_j - \mathbf{w}_l), \quad (18)$$

where  $\kappa$  is simply the Fourier transform of the base measure  $\mu$ , i.e. its characteristic function in the case where  $\mu$  is a probability measure. Hence many standard choices of  $\mu$  lead to tractable normalizing constants. Note that while  $\mathbf{V}$  and  $\kappa$  both may be complex-valued, the normalizing constant as well as the density itself are real-valued.

Various examples of probability measures  $\mu$  and its Fourier transforms that give rise to shift-invariant positive definite kernels  $\kappa$  have been studied in the literature on Random Fourier Features (RFF) [38]. Here too, while the functional form of the expressions is identical, the role between the data and the parameters is reversed: in RFF, one is interested in approximating a given kernel on the data instances, by considering a probability measure on the frequency space which is that kernel's inverse Fourier transform.

**Remark 3.** *If we restrict our attention to real-valued  $\mathbf{V}$  and thus explicitly maintain biases inside the parameterisation, this model can also be realised with stacking cos and sin activation so that they share the readout parameters since*

$$\sum_{r=1}^m \left| \sum_{j=1}^n v_{rj} \exp(i(\mathbf{w}_j^\top \mathbf{x} + b_j)) \right|^2 = \|\mathbf{V} \cos(\mathbf{W}\mathbf{x} + \mathbf{b})\|^2 + \|\mathbf{V} \sin(\mathbf{W}\mathbf{x} + \mathbf{b})\|^2.$$

## G Discrete and mixed continuous/discrete SNEFYs

### G.1 Discrete SNEFYs via series

Closed-form expressions for normalising constants of discrete distributions are only advantageous if their computation costs less than naively summing up unnormalised density over all possible states. This is the case if the support has very large or infinite cardinality. An example that extends the classical Poisson distribution in exponential family form is given in Table 1. Here we discuss some other settings.

**NNK as a convergent series** Suppose the support  $\mathbb{X}$  is discrete and let  $h(x)$  be a nonnegative function corresponding with the discrete base measure  $\mu$ . The NNK is given by

$$k_{\sigma, \mathbf{t}, \mu}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \sum_{\mathbf{x} \in \mathbb{X}} \sigma(\mathbf{w}_1^\top \mathbf{t}(\mathbf{x}) + b_1) \sigma(\mathbf{w}_2^\top \mathbf{t}(\mathbf{x}) + b_2) h(\mathbf{x}).$$

Such NNKs often resemble known convergent series.

**Fourier series** For example, if  $\sigma = \cos$  and  $C_x = h(x)$  are some coefficients of a convergent Fourier series, then

$$k_{\sigma, \mathbf{t}, \mu}(\theta_i, \theta_j) = \frac{1}{2} \sum_{x=0}^{\infty} C_x \left( \cos(t(x)(w_1 - w_2) + (b_1 - b_2)) + \cos(t(x)(w_1 + w_2) + (b_1 + b_2)) \right)$$

is just a series representation of a sum of two periodic functions. For example, if  $C_x = \frac{2(-1)^{x+1}}{\pi x}$  for  $x \geq 1$  and  $C_x = 0$  otherwise, and  $t(x) = 2\pi x$ , the NNK is a sum of two sawtooth waves with frequencies  $w_1 - w_2$  and  $w_1 + w_2$  and phase offsets  $b_1 - b_2$  and  $b_1 + b_2$ . Other examples include rectified sine waves, square waves and triangular waves. This extends to periodic functions with convergent multivariate Fourier series.

## G.2 Mixed continuous/discrete SNEFYs

Mixed distributions can be obtained by choosing the base measure to be a mixed continuous distribution. For example, choose  $\mathbb{X} = \mathbb{R}^d$ ,  $\mu(d\mathbf{x}) = \frac{1}{2}(\Phi(\mathbf{x}) + \delta(\mathbf{x})) d\mathbf{x}$ , and take  $\sigma$  and  $\mathbf{t}$  to be any of the combinations leading to a closed-form NNGPK (for example,  $\mathbf{t}$  as the identity and  $\sigma$  as the error function, Leaky ReLU, GELU, cosine, or snake). Then

$$\begin{aligned} k_{\sigma, \mathbf{t}, \mu}(\theta_i, \theta_j) &= \frac{1}{2} \left( \int_{\mathbb{X}} \sigma(\mathbf{W}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{W}_j^\top \mathbf{t}(\mathbf{x}) + b_j) \delta(\mathbf{x}) d\mathbf{x} + \right. \\ &\quad \left. \int_{\mathbb{X}} \sigma(\mathbf{W}_i^\top \mathbf{t}(\mathbf{x}) + b_i) \sigma(\mathbf{W}_j^\top \mathbf{t}(\mathbf{x}) + b_j) \Phi(\mathbf{x}) d\mathbf{x} \right) \\ &= \frac{1}{2} \left( \sigma(\mathbf{W}_i^\top \mathbf{t}(\mathbf{0}) + b_i) \sigma(\mathbf{W}_j^\top \mathbf{t}(\mathbf{0}) + b_j) + k_{\sigma, \mathbf{t}, \Phi}(\theta_i, \theta_j) \right), \end{aligned}$$

which is a closed form.