
Learning Exponential Families

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently much attention has been paid to implicit probabilistic models – models
2 defined by mapping a simple random variable through a complex transformation,
3 often a deep neural network. These models have been used to great success for
4 variational inference, generation of complex data types, and more. In most all of
5 these settings, the goal has been to find a *particular member* of that model family:
6 optimized parameters index a distribution that is close (via a divergence or clas-
7 sification metric) to a target distribution (such as a posterior or data distribution).
8 Much less attention, however, has been paid to the problem of *learning a model*
9 *itself*. Here we define implicit probabilistic models with specific deep network
10 architecture and optimization procedures in order to learn intractable exponential
11 family models (*not* a single distribution from those models). These exponential
12 families, which are central to some of the most fundamental problems in probabilis-
13 tic inference, are learned accurately, allowing operations like posterior inference
14 to be executed directly and generically by an input choice of natural parameters,
15 rather than performing inference via optimization for each particular realization
16 of a distribution within that model. We demonstrate this ability across a number
17 of non-conjugate exponential families that appear often in the machine learning
18 literature.

1 Introduction

20 Probability models, the fundamental object of Bayesian machine learning, have long challenged
21 researchers with the tradeoff between tractability and expressivity. Though well understood that a
22 model should be chosen to instantiate a set of assumptions and capture existing domain knowledge
23 [1, 2, 3], for many years too-simple models were chosen for their practical advantages (such as
24 conditional conjugacy), which left much to be desired in terms of expressive performance and
25 scalability of these models.

26 More recently the pendulum has swung, via a resurgence in models which map a latent random
27 variable $w \sim q_0$ through a member of a highly expressive function family $\mathcal{G} = \{g_\theta : \theta \in \Theta\}$, the
28 composition resulting in an *implicit probability model* $\mathcal{M} = \{q(g_\theta \circ w) : \theta \in \Theta\}$ (where $q(\cdot)$ is
29 the pushforward density, i.e. the density induced on the image of the random variable w under
30 the function g_θ). Choosing \mathcal{G} to be a parameter-indexed family of neural networks has both a rich
31 history [4, 5], and has recently been used to produce exciting results for density estimation [6, 7, 8],
32 generation of complex data [9], variational inference [10, 11, 12], and more. A noted advantage of
33 these implicit density network models is that in many cases they make minimal assumptions about
34 the data generative (or posterior inference) process. On the other hand, since these models have
35 been chosen to be generic and flexible, they can lack the classic stipulation that a model instantiates
36 existing domain knowledge. The downsides of a too-flexible model with finite data (albeit large)
37 – and the corresponding bias-variance benefit of a restricted model – are textbook knowledge [13,

§7.3], and work on generalization and compressibility in deep networks suggests that this broad class of function families are indeed quite large, perhaps problematically so [14].

Is all the flexibility of an implicit density network model \mathcal{M} always necessary? Consider the case of variational inference, where a generative model $p(z)p_\beta(X|z)$ (latent z , observed data X) is stipulated in the classic sense to embody modeling assumptions (hierarchical model, topic model, Bayesian logistic regression, etc.). When such a model is intractable, it is increasingly common to deploy an implicit “recognition network” model for variational inference [10], which finds a $q_{\theta^*}(z) \in \mathcal{M}$ such that an evidence bound is optimized with respect to the true posterior $p(z|X)$. However, note the widely recognized fact [15] that many such true posteriors $p(z|X)$ belong to models that can be written as exponential families (albeit intractable, due to the choice of sufficient statistics $t(z)$), of the form: $\mathcal{P} = \left\{ \frac{h(z)}{A(\eta)} \exp \{ \eta^\top t(z) \} : \eta \in H \right\}$. Some effort has been made to learn single members of exponential families from their mean parameters [16], but here we are focused on the natural parameterization and the model itself (not simply members thereof).

Should we be able to learn a tractable approximation to this exponential family model, we would in the very least get the bias-variance benefits of an intelligently restricted model space, and at best would get inference “for free” in the sense that we could evaluate approximate posteriors directly without separate optimization for each dataset encountered (a different form of amortized inference [17, 10, 11, 18]). In this paper we aim to learn a restricted model $\mathcal{Q} = \{q(z; \eta : \eta \in H)\}$ that will be a strict subset of the deep implicit model \mathcal{M} and will closely approximate a target exponential family \mathcal{P} . Note the critical difference between this aim and much of the literature that seeks to learn a density $q_\theta^* \in \mathcal{M}$ (we explore this distinction in depth both algorithmically and empirically).

To proceed, we must first specify a set of models $\mathcal{Q} = \{Q_\phi : \phi \in \Phi\}$, from which we can learn a single model Q_{ϕ^*} , and we must second define a sensible parameter space H of each model. To the first, we restrict Θ , the parameter space of \mathcal{M} , to be itself the image of a second deep *parameter network* family $\mathcal{F} = \{f_\phi : \phi \in \Phi\}$, such that $\{f_\phi(\eta) : \eta \in H\} \subset \Theta$. The second part is answered immediately by our choice of target \mathcal{P} , an exponential family which by definition has *natural* parameterization $\eta \in H$. Thus, appealingly, we know that H is precisely the correct parameter space for \mathcal{Q} (as it defines \mathcal{P}), and that the image of H under f_ϕ will be of the correct dimensionality within the codomain Θ ; approximation error between \mathcal{Q} and \mathcal{P} will be caused by the flexibility and learnability of the parameter network f_ϕ and the density network $g_{f_\phi(\eta)}$.

We define this two-network architecture, which we term an *exponential family network* (EFN), and we specify a stochastic optimization procedure over a variant of the typical Kullback-Leibler divergence. We then demonstrate the ability of EFNs to approximately learn exponential families, both known tractable families and well-used intractable families, including hierarchical Dirichlet and truncated normal Poisson families. Finally we demonstrate the utility of this approach in an example inferring the posterior distribution of the latent intensity of a point-process, given neural spike train data. In all, our contributions include:

- a novel implicit model: a two-network deep architecture to learn a probability model along with a doubly stochastic optimization that samples over both natural parameters (the family member to be learned) and data points (observations of the target density);
- analysis of the connections between approximately learning a model and approximate variational inference, and an empirical study that gives insight to possible improvements to variational inference;
- empirical results confirming performance against ground truth in known tractable exponential families and in common intractable exponential families.

2 Exponential family networks

To define exponential family networks (EFN), we begin with relevant context for our modeling choice of exponential families (§2.1). We then describe the primary network architectural constraint and the background we leverage to satisfy that constraint (§2.2). We then introduce EFN in detail, including the optimization algorithm used for learning (§2.3). The similarities with variational inference are then explored in depth in §2.4.



Figure 1: (A) Graphical model for conditionally iid sampling from an exponential family likelihood. (B) Hierarchical Dirichlets – prior $p_0(z)$ (top), three sample conditional Dirichlet datasets X of $N = 2, N = 20, N = 100$ (middle), and three corresponding posteriors that themselves form an exponential family \mathcal{P} (bottom). (C) Architecture for exponential family network (EFN) – density network running top to bottom; parameter network running right to left.

2.1 Exponential families as target model \mathcal{P}

We will focus on a fundamental problem setup in probabilistic inference, that of a latent variable $z \in \mathcal{Z}$ with prior belief $p_0(z)$, and where we observe a dataset $X = \{x_1, \dots, x_N\} \subset \mathcal{X}$ as conditionally independent draws given z . Updating our belief with data produces the posterior $p(z|X) \propto p_0(z) \prod_{i=1}^N p(x_i|z)$. This setup is shown as a graphical model in Figure 1A.

In rare cases these posterior distributions are tractable due to either known conjugacy or to careful historical work (often an inversion, transformation-rejection, or similar custom numerical strategy) that has made these distributions computationally indistinguishable from tractable [19]. It is intriguing then to reflect upon the success that deep networks have offered to function approximation, and ask to what extent we can automate this numerical process, widening the class of effectively tractable distributions.

If we restrict our attention to priors and likelihoods that belong to exponential families $\mathcal{P} = \left\{ \frac{h(\cdot)}{A(\eta)} \exp \{ \eta^\top t(\cdot) \} : \eta \in H \right\}$, the posterior can be also viewed as an exponential family, albeit intractable [15]. For simplicity we will hereafter suppress the base measure $h(\cdot)$. Consider:

$$p_0(z) = \frac{1}{A_0(\alpha)} \exp \{ \alpha^\top t_0(z) \} \quad , \quad p(x_i|z) = \frac{1}{A(z)} \exp \{ \nu(z)^\top t(x_i) \} \quad ,$$

where $t(\cdot)$ is the sufficient statistic vector, and $\nu(z)$ is the natural parameter of the likelihood in natural form [20]. The posterior then has the form:

$$p(z|x_1, \dots, x_N) \propto \exp \left\{ \left[\begin{array}{c} \alpha \\ \sum_i t(x_i) \\ -N \end{array} \right]^\top \left[\begin{array}{c} t_0(z) \\ \nu(z) \\ \log A(z) \end{array} \right] \right\} \quad ,$$

which again is an exponential family, albeit intractable.

To give a concrete example, consider the hierarchical Dirichlet – a Dirichlet prior $z \sim \text{Dir}(\alpha)$ (of dimension $|\mathcal{Z}|$) with conditionally iid Dirichlet draws $x_i|z \sim \text{Dir}(\beta z)$, which has been considered historically [21], and is perhaps most notable for its nonparametric extension [22] (and has relevance for multi-corpus extensions of topic models [23, 24]). Figure 1B shows the prior for

110 a given α (top), and three examples of datasets that could arise via this generative model (mid-
 111 dle). A set of basic manipulations shows the hierarchical Dirichlet posterior $p(z|X)$ to be itself an
 112 exponential family with natural parameter $\eta = [\alpha - 1, \sum_i \log(x_i), -N]^\top$ and sufficient statistic
 113 $t(z) = [\log(z), \beta z, \log(B(\beta z))]^\top$.¹ The corresponding posteriors are shown in Figure 1B (bottom).

114 Note importantly that, because the likelihood was chosen to be an exponential family (which is closed
 115 under sampling), this form will not change for any choice of $|Z|$ -dimensional hierarchical Dirichlet
 116 – any draw from the prior, any N , or any particular realization of observed data X (technically the
 117 prior need not be exponential family, but we leave it as such for simplicity). The exponential family
 118 is clearly sufficient for this property, and the Pitman-Koopman Lemma further clarifies that it is also
 119 necessary (under reasonable conditions) [20, §3.3.3].

120 The critical observation here is that, if we can approximately learn an intractable exponential family
 121 (the model itself), then it becomes trivial to perform posterior inference: we simply use the dataset to
 122 index into the natural parameter η of the intractable family, and the posterior distribution is produced.
 123 This is the goal of EFN.

124 2.2 Density networks as generic approximating family \mathcal{M}

Implicit probability models, which we will use for our approximating model family \mathcal{M} , can be
 defined by any base random variable $w \sim p_0$ mapped through any measurable, parameter-indexed
 function family $\mathcal{G} = \{g_\theta : \theta \in \Theta\}$; we denote the induced density on $z = g_\theta(w)$ as $q_\theta(z)$. Though
 trivial to sample from $q_\theta(z)$ for any choice of family \mathcal{G} , we here additionally require that we be able to
 explicitly calculate $q_\theta(z)$. This goal can be readily achieved by designing \mathcal{G} to contain only bijective
 functions, ideally with a Jacobian form that is convenient to compute. Designing that bijective \mathcal{G} as a
 deep neural network family, as we do here, is a well-established idea that has recently seen many
 variants and applications [5, 25, 26, 7, 6, 27, 28, 8, 29]. Specifically, let $z = g_\theta(w) = g_L \circ \dots \circ g_1(w)$
 for bijective vector-valued functions g_ℓ (where for clarity we have suppressed the dependence of each
 on θ), and denote $J_\theta^\ell(z)$ as the Jacobian of the function g_ℓ at the layer activation corresponding to z .
 Then we have:

$$q_\theta(z) = q_0(g_1^{-1} \circ \dots \circ g_L^{-1}(z)) \prod_{\ell=1}^L \frac{1}{|J_\theta^\ell(z)|}.$$

125 The specific form of the layers g_ℓ can be chosen based on empirical considerations; we clarify our
 126 choice in §3. For the remainder (and to avoid confusion when we introduce a second network) we call
 127 this deep bijective neural architecture the *density network*; this network is shown vertically oriented
 128 (flowing from w down to z) in Figure 1C.

129 This density network induces the model $\mathcal{M} = \{q(g_\theta \circ w) : \theta \in \Theta\}$, which previous work has
 130 searched to find a single optimized distribution (such as a posterior or data generative density), on the
 131 assumption and subsequent empirical evidence that the target exponential family member is close to
 132 (or approximately belongs to) \mathcal{M} . We make the same assumption for the exponential family itself
 133 and seek to intelligently restrict \mathcal{M} in order to learn the exponential family.

134 2.3 Exponential family networks as approximating model \mathcal{Q}

135 Having introduced our target model \mathcal{P} , an exponential family with natural parameters $\eta \in H$, and
 136 the density network family \mathcal{M} , we now seek to learn $\mathcal{Q} \approx \mathcal{P}$, where $\mathcal{Q} \subset \mathcal{M}$. To do so we will
 137 parameterize θ , the parameters of the density network, as the image of a second *parameter network*
 138 family $\mathcal{F} = \{f_\phi : H \rightarrow \Theta, \phi \in \Phi\}$. This network is shown flowing from right to left in Figure 1C.
 139 Using a second meta-network to aid or restrict network learning has been used in a variety of settings;
 140 a few examples include parameterizing the optimization algorithm in the so-called “learning to learn”
 141 setting [30], and a more closely related work that used a second network to condition on observations
 142 for local latent variational inference [27], a connection which we explore closely in the following
 143 section.

144 Any choice of parameter network parameters ϕ induces a $|H|$ -dimensional submanifold (the image
 145 $f_\phi(H)$) of the density network parameter space Θ , and as such defines a restricted model $\mathcal{Q}_\phi =$

¹To be clear this model is an exponential family if β is fixed or treated as a latent variable; this fact however
 will not be important for the development of this paper.

146 $\{q_{f_\phi}(z; \eta) : \eta \in H\} \subset \mathcal{M}$; by our choice of H as the natural parameter space of the exponential
 147 family target \mathcal{P} , this model restriction is at least of the correct dimensionality. Our goal then is to
 148 search over the implied set of models $\mathbb{Q} = \{Q_\phi : \phi \in \Phi\}$ to find an optimal ϕ^* such that $Q_{\phi^*} \approx \mathcal{P}$.

Given the connections between the exponential family and Shannon entropy, we will measure the error between Q_ϕ and \mathcal{P} with Kullback-Leibler divergence. Consider for the moment a fixed choice of natural parameter η ; we seek to minimize, over ϕ :

$$D(q_\phi(z; \eta) || p(z; \eta)) \propto \mathbb{E}_{q_\phi} \left(\log q_\phi(z; \eta) - \eta^\top t(z) \right) = \mathbb{E}_{q_\phi} \left(q_0(g_\theta^{-1}(z)) + \sum_{\ell=1}^L \log |J_\theta^\ell(z)| - \eta^\top t(z) \right),$$

149

150 where again we note that $\theta = f_\phi(\eta)$, and thus for a fixed eta, this objective depends only on ϕ . Indeed,
 151 the target $\eta^\top t(z)$ is linear in η (an obvious restatement of the log-linear exponential family form),
 152 giving us some hope that we may be able to learn this model. As a side note, this objective can also
 153 produce approximations of the log partition (as the intercept term implied by this linear target), which
 154 we have found to be reasonably accurate, though nuanced schemes are likely appropriate [31]; we do
 155 not explore that further here.

156 Of course we seek to approximate not just a single target exponential family member ($p(z; \eta)$ for
 157 a fixed η), but rather the entire model $\mathcal{P} = \{p(z; \eta) : \eta \in H\}$. For optimization we thus need to
 158 introduce a distribution $p(\eta)$ (for sampling), leading to the objective:

$$\operatorname{argmin}_{\phi} \mathbb{E}_{p(\eta)} (D(q_\phi(z; \eta) || p(z; \eta))) = \operatorname{argmin}_{\phi} D(q_\phi(z; \eta)p(\eta) || p(z; \eta)p(\eta)).$$

159 Unbiased estimates of this objective are immediate: $q_\phi(z; \eta)$ is sampled by computing calculating
 160 the density network parameters $\theta = f_\phi(\eta)$ (using the parameter network), sampling the latent
 161 $w \sim p_0(w)$, and running that w through the density network; $p(\eta)$ is user defined and thus trivial to
 162 sample. Stochastic optimization can then be carried out on the estimator:

$$\mathbb{L}(\phi) = \frac{1}{K} \frac{1}{M} \sum_{k=1}^K \sum_{m=1}^M \left(q_0(g_{\theta^k}^{-1}(z^m)) + \sum_{\ell=1}^L \log |J_{\theta^k}^\ell(z^m)| - \eta_k^\top t(z^m) \right),$$

163 where $\theta^k = f_\phi(\eta_k)$. Successful optimization over ϕ should thus result in $Q_{\phi^*} \in \mathbb{Q}$ that accurately
 164 approximates the target exponential family; that is, $Q \approx \mathcal{P}$. We call this two-network architecture
 165 and optimization an exponential family network (EFN). What remains for empirical implementation
 166 is to make particular choices of hyperparameters, network layers, and optimization algorithm, which
 167 we specify in §3 below.

168 2.4 Relation to variational inference

169 We have already covered related work; here we scrutinize EFNs in terms of VI.

170 We are interested in perhaps the most classic inference problem:

$$p(z|x) \propto p(z) \prod_{i=1}^n p(x_i|z)$$

171 shown with the attached plate model (not local latents). Supposing as is often the case that the
 172 likelihood is a member of the s exp fam, we have:

$$p(z|x) \propto \exp \left\{ \left[\sum_{i=1}^n s(x_i) \right]^\top [t(z)] + g_0(\alpha, z) \right\}$$

173 Important to distinguish carefully from VI. In a sense VI does parameterize a family: given data,
 174 you get local variational parameters and that parameterizes a density (like a regular VAE). Inference
 175 networks are exclusively used to data to amortize with a global set of parameters a variational

distribution, not a model. Of course it is in a sense a model, but that's a bunch of normals. The sampling mechanism is easy (Gaussian).

where the natural parameters of the sampling distribution are indexed by the latent parameter on which we want to inference (z). Here I've written the prior as arbitrary, and possibly not exp fam, which is fine, since this is still an exp fam in the sense of, for a fixed α , the function g_0 can just be viewed as a sufficient statistic. Even if α is not fixed though, we can sample over that too to learn the whole fam (but maybe not if we want to infer it?). Regardless, life is simpler to make sense of if we take an exp fam prior $g_0(\alpha, z) = \alpha^\top t_0(z)$, and then the desired posterior is an intractable exp fam, but still just an exp fam.

Note: consider changing all z to θ to remind the average reader that we're doing real bayesian inference and not just run of the mill VI with local latents in a nonlinear dimension reduction setting. Perhaps an important reminder that most all of VAE and such are for inference of local latents, and that's a little bit too bad. We fix that.

Another key idea that EFNs enable is to ask if learning the $\theta(\eta)$ network leads to better VI in terms of inference networks, since it is apparently appropriately regularized and can just take suff stats. That's testable if we have time.

In a restricted technical sense, rather close: VAE and other black box VI that uses reparameterization results in a conditional density $q_\phi(z|x)$. If we consider η as x , then sure yes the previous stuff specifies a model $\mathcal{Q}_{VAE} = \{q_\phi(z|x) : x \in X\}$. But that's a little silly, and any way that is very often a normal family with variational parameters specified by (a deep function of) x . Much closer is Figure 2 in Rezende and Mohamed, where like here they use a network to index the *parameters* of the normalizing flow. In that case it's a function of x the observation, and as such that network is an inference network; here it's a function of η and as such is a parameter network. That's just nomenclature, so naturally the next question is do they differ at some other level. Yes, distinctly. The other term implied in a VI (or norm flow VAE style as they use) is the expected log joint $E_{q_{\phi(x)}}(\log p_\theta(x, z))$. Now sure that's a loss function on x, z , so then when we look at that same term in EFN we see $E_{q_{\phi(\eta)}}(\eta^\top t(z))$, which sure also looks like a loss function on η, z . And yes, they are both unnormalized (in the sense that VI is an ELBO / joint $p(x, z)$ and EFN lacks the normalizer because it's constant, so we're not getting a KL estimate). A picky difference is that the exp family doesn't really correspond to a proper unnormalized log joint (though I suppose it could), as there is not a prior on η in the objective (but is that just ignoring $p(\eta)$ in our sampling scheme?). But yes if we want to be reductionist and pedantic [use nicer words] in general we could see this as a specific case where $x = \eta$ and thus we are learning a family just as in the inference case. Or rather, we are putting the data in as sufficient stat (computation of natural parameters), but that's nonobvious. And for example we are giving in the bayesian logistic regression example full datasets for inference instead of single data points. To make this as close as possible, we write $p(\eta|z) = \frac{1}{A(t(z))} \exp \{\eta^\top t(z)\}$. That's the "likelihood" of an EFN in some wonky sense. So this reveals the mechanical differences: first, $t(z)$ is not a deep generative model with parameters θ , but rather it is a fixed set of sufficient statistics that define the exp fam. Next, there is no clear prior $p(z)$, which is critical to understanding how VI behaves (see Hoffman and Johnson ELBO surgery paper, also Duvenaud's <https://arxiv.org/pdf/1801.03558.pdf>). So yes there is a hand wavy sense in which EFN is a specific case of norm flow, but of course it is. And anyway norm flow is a specific case of a DNN architecture or Helmholtz machine or deep density network (Ripple and Adams). This is just rambling but good to have all perspective here. Ok so what to do? First, then we need to produce really compelling results focusing on when learning an exp fam is key. Second we need some very tight language to draw this distinction without seeming a small tweak on normalizing flows. One way to do this is the restricted model class argument, a la Fig 7.2 in Hastie and Tibshirani. Another is to actually produce a conditional exp fam, as in something indexed on both x and η . Third, possible novelties in norm flows, like triple spinners or other better choices than planar flows (yuck).

Another point is that it's unknown if posterior contraction can be well modeled. As in, we know that most VI NF type things are conditioned on a single data point, so the posterior variance can tend to be rather homogenous. One more contribution is to offer that contraction study; as we get more data points we will get more posterior contraction, so this tests the ability of this model to learn that.

Key distinctions:

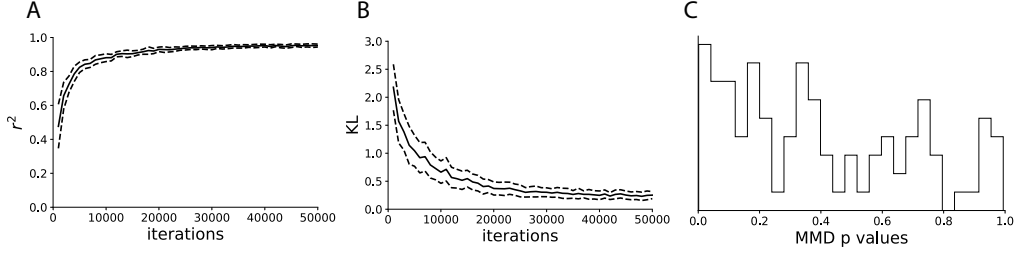


Figure 2: 25-dimensional Dirichlet exponential family network. A.) Distribution of r^2 between the sufficient statistics and log-probability across choices of η throughout optimization. B.) Distribution of KL divergence across choices of η throughout optimization. C.) Distribution of maximum mean discrepancy p-values between EFN samples and ground truth after optimization [38].

230 • narrow mechanical sense this is VI with an observation of the natural parameters, namely
 231 the sample exp fam over all data. but that’s pedantic.

232 • no generative model in the usual sense: yes, we can consider a prior and then some
 233 observation model as the genrative model, but in any event it’s not a neural net.

234 • we lack a finite data set X , so the objective is technically different. We stipulate a distribution
 235 and then this is expectation over that model space, a KL or a KL to the broader joint with η .
 236 This is concretely different, as we typically use a fixed size dataset X so we can calculate
 237 the ELBO over the

238 Latest key distinctions:

239 • prior is in parameter network, unlike essentially all others, even if you take a narrow view
 240 that $\sum_i t(x_i)$ is a single data point. Prior has been recognized for mattering in the ELBO,
 241 though this sentence is a dubious distinction [32, 33] (dubious need for these refs)

242 • data is given via an assumption of sufficiency, namely in natural form [20], not in x form. Of
 243 course this is sensible as in some settings we don’t know the natural form of the generative
 244 model, but that’s a key difference with SVI; plenty of those models are not deep nets (and
 245 shouldn’t be, if there is an intent of statistical inference rather than nonlinear dimensionality
 246 reduction / autoencoding) and there we *do* know the natural parameters.

247 More generally there has been a lot of attention to making these more flexible in structured variational
 248 inference. Any generalization of this is also dandy though, so could use a mean field approach
 249 (standard) or any of the things that go beyond mean field, either classically [34, 35]; this is called
 250 structured variational inference or newer stuff [36] [37], to name but a few.

251 3 Results

252 Introductory remarks and then comments about architectural particulars, including planar flow
 253 networks of [27]. Note Number of panar flows is always D (intrinsic dimensionality of flows), units
 254 per layer ramping is always the same function of D . The number of layers in the theta network is
 255 always a function of D - will probably just always use 8 layers. Remember

256 NF1: do full norm flow “variational inference” (explore all of ϕ space with the full flow network
 257 model \mathcal{Q}), which is to say $\arg \min_{\phi} KL(q_{\phi}||p)$.

258 EFN1: be literal to Figure 1C, give the sufficient statistics of that $K=1$ dataset, and learn an EFN from
 259 scratch. This alternative is important because it is the most specific (but kind of annoying, hence
 260 alternative 1) interpretation of norm flow VI paper.

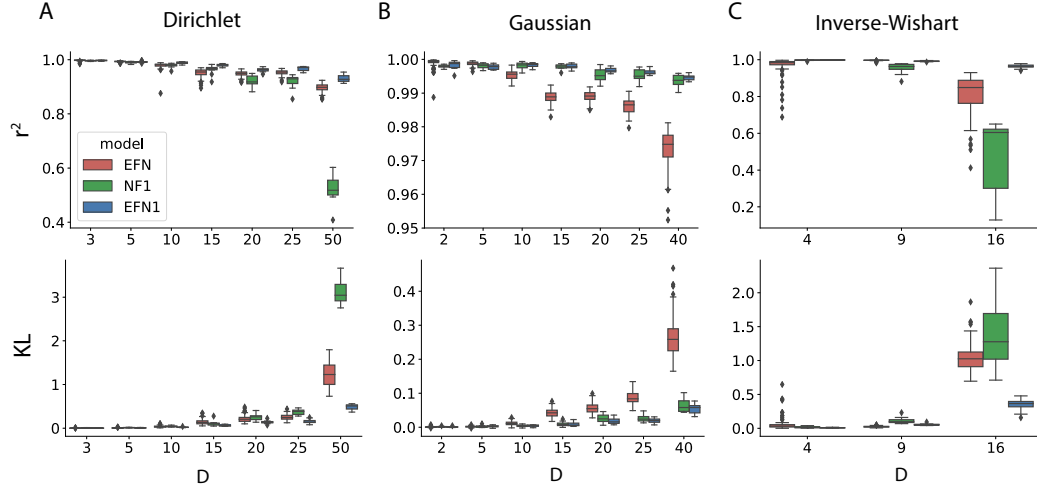


Figure 3: Scaling exponential family networks. A.) Dirichlet. B.) Gaussian C.) Inverse-Wishart

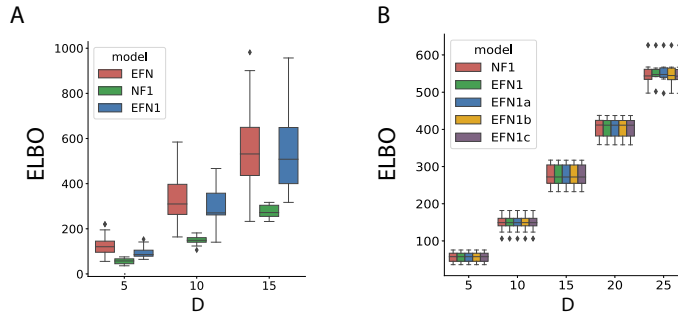


Figure 4: Scaling Dir-Dir

3.1 Tractable exponential families

3.2 Intractable exponential families

Hierarchical Dirichlets Hierarchical dirichlets are useful and have some history; most notable is with the Hierarchical Dirichlet Process [22], but historically this was done in the finite case also [21]. Here is some math. Note that this matters for multi-corpus LDA generally as well [23, 24].

Truncated- and log-normal Poisson used a lot [39][40][41, 42]

Figure 4:

EFN in intractable exp fams (connecting to above, but with hard distribs and the ELBO)

Panel A: Dir-Dir ELBO by dimensionality for NF1 and EFN and EFN1

Panel B: Dir-Dir ELBO by dimensionality for EFN1 vs EFN1a vs 1b vs 1c vs NF1 (with $N = 1$ data point)

273 3.3 Neural spike train analysis

274 Figure 5:

275 Panel A TNP picture example of prior and posterior with a few samples, just for feel good

276 PANEL B: ELBO on held out data as a function of R , for a middle choice of training dataset size N
277 and D .

278 PANEL C: ELBO on held out data as a function of N , for a middle choice of number of samples in
279 the posterior R .

280

281 PANEL D (optional): (ELBO EFN - ELBO NF1) as a surface plot as a function of R, N . That is,
282 positive places is where EFN outperforms, negative NF1.

283 The key point with these is that, while you have the *same exact* flow network architecture, now you
284 have to optimize over ϕ with a limited single dataset. Learning a restricted model space is good for
285 the bias-variance tradeoff! Do this many times so that variance will become clear.

286 —other thoughts— Real data analysis and posterior inference. **Key real data result on TNP.**

287 Get some data from CRCNS that has many spike trains x_i for $i = 1, \dots, N$ (ask Gabriel, as he has
288 done some poking around recently; or look at some of the above TNP/LNP refs).

289 Those spike trains should be conditionally independent draws from the same underlying intensity
290 function z . (for example, trials under the same stimulus)

291 Bin the length of time T into $\approx 20 - 30$ equally spaced time bins. Thus z is now a vector in \mathbb{R}^{20} .

292 Now each spike train x_i is a conditionally independent Poisson vector observation, with rate vector z .

293 Learn the 20 dimensional TNP exp fam, without any regard to this dataset X .

294 No: Panel No: TNP ELBO by dimensionality for NF1 and EFN and EFN1

295 Panel A TNP picture example of prior and posterior with a few samples, just for feel good

296

297 **Now we want to learn the posterior $p(z | \text{some fixed number } R \text{ of data points})$.**

298 To do this for an EFN, just plug in those R points x_{i_1}, \dots, x_{i_R} and the prior as a natural parameter,
299 and job done.

300 To do this for an NF1, train a VI model by taking the log joint with R data points, then go through
301 and resample R points every time from your training dataset with N data points.

302 **PANEL A: ELBO on held out data as a function of R , for a middle choice of training dataset
303 size N .**

304 **PANEL B: ELBO on held out data as a function of N , for a middle choice of number of
305 samples in the posterior R .**

306 **PANEL C: (ELBO EFN - ELBO NF1) as a surface plot as a function of R, N . That is, positive
307 places is where EFN outperforms, negative NF1.**

308 The key point with these is that, while you have the *same exact* flow network architecture, now you
309 have to optimize over ϕ with a limited single dataset. Learning a restricted model space is good
310 for the bias-variance tradeoff! Do this many times so that variance will become clear. **Panel C v2:**

311 **Possibly want to explicitly plot variance of EFN and NF1 to focus on the variance tradeoff**

312 **Panel C v3: change time bin granularity from 10 to 50 to show how this story changes in D .
313 My thought is that all will be exhausted by dimensionality sweeps by this point, so no.**

314 also Notice one pain here is that these panels requires training a new EFN1 at every choice of N and
315 R (but only one EFN). Sorry.

316

317 We hope and expect this will show that when the dataset gets small, this "traditional VI" will get
318 arbitrarily bad (can't learn a network); eventually, there will be so much data that the VI will match
319 or outperform the EFN... outperform because VI can focus specifically on this distribution rather than
320 over the whole family, so the EFN has less effective data for this η (but not because it has a broader
321 range of models, since we believe the EFN contains the closest member). Performance metric should
322 be ELBO on some held out data or something like that (it's a posterior, so log likelihood doesn't
323 really make sense). Test data anyway. Check VI papers for usual metrics. A key point to make
324 here is that one great virtue of EFNs is is learning a restricted model, which should demonstrate the
325 usual bias-variance tradeoff (see for example Hastie and Tibshirani book, Fig 7.2). Or Figure 4 is
326 bias-variance and some sample posteriors in 2-d (showing how nicely it works), and then Fig 5 is the
327 above performance, with both train and test.

328 This will be for one real example X . As such, to get error bars, just take a big dataset and randomly
329 subsample the test set. Then the posterior performance is really for that very dataset, so the sem is
330 coherent and the right thing to calculate/show. Important to clarify that doing so *does not* test how
331 well this does across the entire exp fam, but just this one posterior. ((To test that, we would do it in
332 simulation: generate *many datasets* X , then do the above for every one of them. Same computation
333 for EFN (since its just plugging in a dataset), but VI alternatives 1 and 2 now need to be rerun for
334 every dataset. And it's still simulated data, not really offering something fundamentally more than Fig
335 3 (well ok it's an intractable model, but I'm not sure that offers so much)...let's skip that altogether)).

336 **4 Conclusion**

337 Snappy closing remarks!

References

- [1] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*, volume 2. CRC press Boca Raton, FL, 2014.
- [2] Joshua B Tenenbaum, Thomas L Griffiths, and Charles Kemp. Theory-based bayesian models of inductive learning and reasoning. *Trends in cognitive sciences*, 10(7):309–318, 2006.
- [3] Peter McCullagh. What is a statistical model? *The Annals of Statistics*, 30(5):1225–1267, 2002.
- [4] Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.
- [5] D. J. C. MacKay and M. N. Gibbs. Density networks. In *Statistics and Neural Networks*, pages 129–146. Oxford, 1997.
- [6] Benigno Uribe, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems*, pages 2175–2183, 2013.
- [7] Oren Rippel and Ryan Prescott Adams. High-dimensional probability estimation with deep density models. *arXiv preprint arXiv:1302.5125*, 2013.
- [8] George Papamakarios, Iain Murray, and Theo Pavlakou. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2335–2344, 2017.
- [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- [10] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv*, 12 2013.
- [11] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- [12] Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational bayes for non-conjugate inference. In *International Conference on Machine Learning*, pages 1971–1979, 2014.
- [13] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [14] Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P Adams, and Peter Orbanz. Compressibility and generalization in large-scale deep learning. *arXiv preprint arXiv:1804.05862*, 2018.
- [15] 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.
- [16] Gabriel Loaiza-Ganem, Yuanjun Gao, and John P Cunningham. Maximum entropy flow networks. *International Conference on Learning Representations*, 2017.
- [17] Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 36, 2014.
- [18] Andreas Stuhlmüller, Jacob Taylor, and Noah Goodman. Learning stochastic inverses. In *Advances in neural information processing systems*, pages 3048–3056, 2013.
- [19] Luc Devroye. *Non-uniform random variate generation*. Springer-Verlag, New York, 1986.
- [20] Christian Robert. *The Bayesian choice: from decision-theoretic foundations to computational implementation*. Springer Science & Business Media, 2007.
- [21] David JC MacKay and Linda C Bauman Peto. A hierarchical dirichlet language model. *Natural language engineering*, 1(3):289–308, 1995.

- [22] Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Hierarchical dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [23] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- [24] Jonathan K Pritchard, Matthew Stephens, and Peter Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000.
- [25] Leemon Baird, David Smalenberger, and Shawn Ingkiriwang. One-step neural network inversion with pdf learning and emulation. In *Neural Networks, 2005. IJCNN’05. Proceedings. 2005 IEEE International Joint Conference on*, volume 2, pages 966–971. IEEE, 2005.
- [26] Esteban G Tabak, Eric Vanden-Eijnden, et al. Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences*, 8(1):217–233, 2010.
- [27] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.
- [28] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*, 2016.
- [29] Jörn-Henrik Jacobsen, Arnold Smeulders, and Edouard Oyallon. i-revnet: Deep invertible networks. *arXiv preprint arXiv:1802.07088*, 2018.
- [30] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*, pages 3981–3989, 2016.
- [31] George Papamakarios and Iain Murray. Distilling intractable generative models. In *Probabilistic Integration Workshop at Neural Information Processing Systems*, 2015.
- [32] Matthew D Hoffman and Matthew J Johnson. Elbo surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference, NIPS*, 2016.
- [33] Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders. *arXiv preprint arXiv:1801.03558*, 2018.
- [34] Lawrence K Saul and Michael I Jordan. Exploiting tractable substructures in intractable networks. In *Advances in neural information processing systems*, pages 486–492, 1996.
- [35] David Barber and Wim Wiegerinck. Tractable variational structures for approximating graphical models. In *Advances in Neural Information Processing Systems*, pages 183–189, 1999.
- [36] Matthew Hoffman and David Blei. Stochastic structured variational inference. In *Artificial Intelligence and Statistics*, pages 361–369, 2015.
- [37] Dustin Tran, David Blei, and Edo M Airoldi. Copula variational inference. In *Advances in Neural Information Processing Systems*, pages 3564–3572, 2015.
- [38] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- [39] Yuanjun Gao, Evan W Archer, Liam Paninski, and John P Cunningham. Linear dynamical neural population models through nonlinear embeddings. In *Advances in Neural Information Processing Systems*, pages 163–171, 2016.
- [40] Ryan Prescott Adams, Iain Murray, and David JC MacKay. Tractable nonparametric bayesian inference in poisson processes with gaussian process intensities. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 9–16. ACM, 2009.

- 427 [41] John P Cunningham, Krishna V Shenoy, and Maneesh Sahani. Fast gaussian process methods
428 for point process intensity estimation. In *Proceedings of the 25th international conference on*
429 *Machine learning*, pages 192–199. ACM, 2008.
- 430 [42] John P Cunningham, M Yu Byron, Krishna V Shenoy, and Maneesh Sahani. Inferring neural
431 firing rates from spike trains using gaussian processes. In *Advances in neural information*
432 *processing systems*, pages 329–336, 2008.