# **Learning Exponential Families**

#### **Anonymous Author(s)**

Affiliation Address email

#### **Abstract**

1	[SLOFF I NOTES STAGE JUST TO GET THOUGHTS DOWN]
2	Recently much attention has been paid to probabilistic models defined by a deep
3	neural network transformation of a simpler random variable; these implicit genera-
4	tive models have been used to great success across variational inference, generative
5	modeling of complex data types, and more. In essentially all of these settings, the
6	model is specified by the network architecture, and a particular member of that
7	model is chosen to minimize some loss (be it adversarial or information divergence)
8	We treat the problem of learning an exponential family – the model itself, rather
9	than the typical setting of learning a particular member of that model.
10	Many intractable distributions encountered in machine learning belong to expo-
11	nential families. In rare cases these distributions are tractable due to either known
12	conjugacy in the problem setup (such as the normal-inverse-Wishart), or due to care-
13	ful numerical work historically that has made these distributions computationally

ICLODDY MOTES STACE HIST TO SET THOUSELITS DOWN!

## 5 1 Introduction

14

16 People use lots of implicit generative models:

Across machine learning, including ABC [?], GANs [1], VAEs [2, 3], and their many follow-ons (too numerous to cite in any detail), models that specify a distribution via the nonlinear transformation of latent random variable. We prefer and use the terminology of [4], calling such a distribution an *implicit generative model*, defined as:

#### something like eq 1 and 2 in Mohamed: 2016aa, defining $q_{\theta}(z)$

- Also use the proper notation of the density implied by the pushforward measure of the function  $f_{\theta\sharp}$  if useful. Also reference to this being super standard and widespread [5]. The two central uses are at present generative distributions of interesting data types (as in GANs), and for variational inference Regardless, all of these use cases specify a *model* (or variational family)  $\mathcal{Q} = \{q_{\theta}: \theta \in \Theta\}$ , and then minimize a suitable loss  $\mathcal{L}(q,p)$  over  $q \in \mathcal{Q}$ . In the case of VI p is the posterior (or the unnormalized log joint) and  $\mathcal{L}$  is the KL divergence (or so called ELBO), in GAN p is the sample density of a (large) dataset and  $\mathcal{L}$  is the adversarial objective whose details do not matter here.
- 24 All these learn a single member of a family
- Inherent in all the above approaches is an algorithmic procedure to select a *single* distribution  $q_{\theta}(z)$
- from among the *model Q*. Implicit in this effort is the belief that Q is suitably general to contain the
- 27 true distribution of interest, or at least an adequately close approximation.

indistinguishable from tractable (eg the Dirichlet).

- 28 Here we learn the family
- 29 We leverage the natural parameterization of exponential families to derive a novel objective that is
- 30 amenable to stochastic optimization.

- 31 A note on amortization
- 32 Several have pointed out that these IGMs are in fact strictly less expressive than a mean field, at
- least in the conventional VI setting. See for example http://dustintran.com/blog/variational-auto-
- encoders-do-not-train-complex-generative-models (here I like the line "The neural network used in
- 35 the encoder (variational distribution) does not lead to any richer approximating distribution. It is a
- way to amortize inference such that the number of parameters does not grow with the size of the data
- 37 (an incredible feat, but not one for expressivity!) (Stuhlmuller et al., 2013)"). You have to optimize
- for every data point individually, or instead you get to do so in aggregate once in advance (at a much
- 39 higher cost) and then recover that cost over future data points within that distribution (and hence the
- 40 term amortization, though perhaps there is shared statistical power as well) Etc etc what we are doing
- 41 here is amortized amortized inference, in the sense that we are amortizing not the data points, but the
- 42 distribution itself.
- REparameterization trick (Kingma and Welling (2013), Rezende et al. (2014) and Titsias and
- 44 Lazaro-Gredilla 2014).. See also Archer 2015 / Gao 2016 for clean explanation.
- 45 Key for obvious norm flow connection but also a good bibliography and some good historical views
- to Dayan and Gershman and other people who did norm flows. https://arxiv.org/pdf/1505.05770.pdf
- 47 Our contributions include:
- 48 ...
- 49 This should not be confused with "Learning to learn by gradient descent by gradient descent"
- 50 (Andrycowicz et. al. 2016) and similar works.
- 51 ...
- 52 Our results demonstrate
- 53 ...

# 4 2 Learning exponential families

- 55 Why this is important
- Exp fams are awesome and fundamental []. Also [?] rightly point out that many many inference
- 57 problems can be cast as exponential families. Can we cast the VAE encoder network as a suitable
- 58 exp fam... sure I think that's right; the network parameters of z form the statistics, and then the
- observations are eta's.
- 60 Why this is coherent
- $\Theta$  defines quite a big Q, and indeed the subject of compressibility, generalization, etc is of keen
- 62 interest to many [?]. So actually the space of distributions is quite large, and in many cases certainly
- 63 larger than it needs be. Why? Well, we know precisely the parameter space of the exponential family;
- it is defined by the *natural* parameters  $\eta \in \mathbb{R}^p$  (or whatever we choose there).
- 65 Figure 1
- Figure of model space. Yeah that's good. Then graphical model. Note that perhaps Q is too big,
- and a simpler model space (the  $\|\eta\|$  dimensional subspace of  $\Theta$ ) would be better for the usual
- 68 robustness/generalization reasons.
- 69 Aside
- 70 A neat idea is to ask if learning the  $\theta(\eta)$  network leads to better VI in terms of inference networks,
- 71 since it is apparently appropriately regularized and can just take suff stats. That's testable if we have
- 72 time.
- 73 Why Flow Networks
- 74 We choose flow networks [] and [] because duh. That makes these comparisons legitimate and apples
- 75 to apples. Any generalization of this is also dandy though, so could use a mean field approach
- 76 (standard) or any of the things that go beyond mean field, either classically (Saul and Jordan, 1996;
- 77 Barber and Wiegerinck, 1999); this is called structured variational inference. Another way to expand

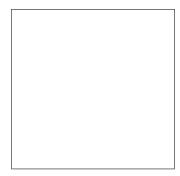


Figure 1: Figure 1: possibly Fig 7.2 bias-variance tradeoff and then benefit of a restricted model from Hastie Book, or similar from W4400 (ask PO for batman permission).

Table 1: Sample table title

	Part	
Name	Description	Size $(\mu m)$
Dendrite Axon Soma	Input terminal Output terminal Cell body	$\begin{array}{c} \sim \! 100 \\ \sim \! 10 \\ \text{up to } 10^6 \end{array}$

- the family is to consider mixtures of variational densities, i.e., additional latent variables within the variational family (Bishop et al., 1998). or newer stuff [] [Tran Copula VI, Hoffman and Blei 2015].
- 80 Couch this in terms of normalizing flows though point out this is not strictly necessary. Note in
- particular Tabak, E. G. and Turner, C. V. A family of nonparametric density estimation algorithms. Communications on Pure and Applied Mathematics, 66(2):145?164, 2013. Tabak, E. G and Vanden-
- 83 Eijnden, E. Density estimation by dual ascent of the log-likelihood. Communications in Mathematical
- Edition F. Density estimation by dual ascent of the log-incliniod. Communications in Mathematical
- Sciences, 8(1):217?233, 2010. A nice line from Rezende and Mohamed is: Thus, an ideal family of
- $^{85}$  variational distributions q(z|x) is one that is highly flexible, preferably flexible enough to contain the
- 86 true posterior as one solution. One path towards this ideal is based on the principle of nor-malizing
- 87 flows (Tabak Turner, 2013; Tabak VandenEijnden, 2010).
- 88 In many situations, statistical inference attempts to learn, at least approximately, a member of an
- 89 exponential family. We often consider this exponential family intractable in the sense that we don't
- 90 know how to normalize or sample from it. Approximate inference, such as variational

## 91 **3 To Do**

#### 92 3.1 SRB

93

97

- set up submission at https://cmt.research.microsoft.com/NIPS2018/
- review and conform to style requirements (see website with template); 8 pages not including refs and acks and appendices.

#### 96 3.2 JPC

- Outline
- Write
- 99 \usepackage[pdftex]{graphicx} ...
- includegraphics[width=0.8\linewidth] {myfile.pdf}

#### 101 Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

#### 04 References

#### 105 References

- 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.
- 110 [2] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 12 2013.
- 111 [3] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- 113 [4] Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. 10 2016.
- 115 [5] Luc Devroye. Non-uniform random variate generation. Springer-Verlag, New York, 1986.
- 116 Stuff on wake sleep and the Helmholtz machine
- 117 Stuff on sampling from Gibbs distributions (max ent models), and sampling from exp fams generally,
- 118 with MCMC and such.
- 119 Flow networks
- 120 Devroye's book.
- Hoffman et al 2013 SVI
- 122 From Blei review on VI. ThedevelopmentofvariationaltechniquesforBayesian inference followed
- two parallel, yet separate, tracks. Peterson and Anderson (1987) is arguably the first variational
- procedure for a particular model: a neural network. This paper, along with insights from statistical
- mechanics (Parisi, 1988), led to a flurry of variational inference procedures for a wide class of models
- (Saul et al., 1996; Jaakkola and Jordan, 1996, 1997; Ghahramani and Jordan, 1997; Jordan et al.,
- 1999). In parallel, Hinton and Van Camp (1993) proposed a variational algorithm for a similar neural
- network model. Neal and Hinton (1999) (first published in 1993) made important connections to the
- expectation maximization (EM) algorithm (Dempster et al., 1977), which then led to a variety of variational inference algorithms for other types of models (Waterhouse et al., 1996; MacKay, 1997).
- Salimans, T. and Knowles, D. (2014). On using control variates with stochastic approximation for
- variational Bayes. arXiv preprint arXiv:1401.1022.
- Salimans, T., Kingma, D., and Welling, M. (2015). Markov chain Monte Carlo and variational
- inference: Bridging the gap. In International Conference on Machine Learning, pages 1218? 1226.
- Ranganath, R., Gerrish, S., and Blei, D. (2014). Black box variational inference. In Artificial
- 136 Intelligence and Statistics.
- 137 Hoffman, M. D. and Blei, D. M. (2015). Structured stochastic variational inference. In Artificial
- 138 Intelligence and Statistics.
- 139 Possibly some of Burda, Y., Grosse, R., & Salakhutdinov, R. (2016). Importance Weighted Autoen-
- coders. In International Conference on Learning Representations. Damianou, A. C., & Lawrence, N.
- 141 D. (2013). Deep Gaussian Processes. In Artificial Intelligence and Statistics. Dayan, P., Hinton, G. E.,
- Neal, R. M., & Zemel, R. S. (1995). The Helmholtz Machine. Neural Computation, 7(5), 889?904.
- http://doi.org/10.1162/neco.1995.7.5.889 Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density
- estimation using Real NVP. arXiv.org. Harville, D. A (1977). Maximum likelihood approaches
- 145 to variance component estimation and to related problems. Journal of the American Statistical
- Association, 72(358):320?338. Hinton, G. and Van Camp, D (1993). Keeping the neural networks
- 47 simple by minimizing the description length of the weights. In Computational Learning Theory,

pp. 5?13. ACM. Johnson, M. J., Duvenaud, D., Wiltschko, A. B., Datta, S. R., & Adams, R. P. (2016). Composing graphical models with neural networks for structured representations and fast 149 inference. arXiv.org. Kingma, D. P., & Welling, M. (2014). Auto-Encoding Variational Bayes. In 150 International Conference on Learning Representations. Kingma, D. P., Salimans, T., & Welling, M. 151 (2016). Improving Variational Inference with Inverse Autoregressive Flow. arXiv.org. Louizos, C., 152 & Welling, M. (2016). Structured and Efficient Variational Deep Learning with Matrix Gaussian 153 Posteriors. In International Conference on Machine Learning. Maaloe, L., Sonderby, C. K., Sonderby, 154 S. K., & Winther, O. (2016). Auxiliary Deep Generative Models. In International Conference 155 on Machine Learning. MacKay, D. J., & Gibbs, M. N. (1999). Density networks. Statistics and 156 neural networks: advances at the interface. Oxford University Press, Oxford, 129-144. Mnih, A., & 157 Rezende, D. J. (2016). Variational inference for Monte Carlo objectives. In International Conference 158 on Machine Learning. Ranganath, R., Tran, D., & Blei, D. M. (2016). Hierarchical Variational 159 Models. In International Conference on Machine Learning. Rezende, D. J., Mohamed, S., & Wierstra, 160 D. (2014). Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In 161 International Conference on Machine Learning. Salimans, T., Kingma, D. P., & Welling, M. (2015). 162 Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. In International Conference 163 on Machine Learning. Salakhutdinov, R., Tenenbaum, J. B., and Torralba, A (2013). Learning with 164 hierarchical-deep models. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 35 165 (8):1958?1971. Stuhlmuller, A., Taylor, J., & Goodman, N. (2013). Learning Stochastic Inverses. 166 In Neural Information Processing Systems. Tran, D., Blei, D. M., & Airoldi, E. M. (2015). Copula 167 variational inference. In Neural Information Processing Systems. Tran, D., Ranganath, R., & Blei, D. 168 M. (2016). The Variational Gaussian Process. International Conference on Learning Representations. 169 Waterhouse, S., MacKay, D., and Robinson, T (1996). Bayesian methods for mixtures of experts. In 170 Neural Information Processing Systems. 171