
Bayesian Inference in Augmented Bow Tie Networks

Jimmy T.H. Smith
Stanford University
jsmith14@stanford.edu

Dieterich Lawson
Stanford University
jdlawson@stanford.edu

Scott W. Linderman
Stanford University
swl1@stanford.edu

Abstract

We develop a deep generative model that generalizes feed-forward, rectified linear neural networks with stochastic activations. We call these models *bow tie networks* because of the shape of their activation distributions. We leverage the Pólya-gamma augmentation scheme to render the model conditionally conjugate and derive a block Gibbs sampling algorithm to approximate the posterior distribution over activations and model parameters. The resulting algorithm is massively parallelizable. We show a proof-of-concept of this model and Bayesian inference algorithm on a variety of standard regression benchmarks.

1 Introduction

Consider a deep generative model for nonlinear regression. Let $\mathbf{x}_n \in \mathbb{R}^{D_x}$ denote the inputs, $\mathbf{y}_n \in \mathbb{R}^{D_y}$ the outputs, and $\mathbf{a}_n = \{\mathbf{a}_{n,l}\}_{l=1}^L$ with $\mathbf{a}_{n,l} \in \mathbb{R}^{D_l}$ the latent activations at each of L intermediate layers. We model the joint distribution as,

$$p(\mathbf{y}_n, \mathbf{a}_n \mid \mathbf{x}_n, \boldsymbol{\theta}) = \left[\prod_{l=1}^L \mathcal{N}(\mathbf{a}_{n,l} \mid \boldsymbol{\mu}_{n,l}, \boldsymbol{\Sigma}_l) \right] \mathcal{N}(\mathbf{y}_n \mid \mathbf{u}_{n,L+1}, \boldsymbol{\Sigma}_{L+1}) \quad (1)$$

$$\boldsymbol{\mu}_{n,l} \triangleq f(\mathbf{u}_{n,l}) \quad \mathbf{u}_{n,1} \triangleq \mathbf{W}_1 \mathbf{x}_n + \mathbf{b}_1, \quad \mathbf{u}_{n,l} \triangleq \mathbf{W}_l \mathbf{a}_{n,l-1} + \mathbf{b}_l \text{ for } l > 1, \quad (2)$$

where f is a nonlinear function applied element-wise. The parameters consist of the weights, biases, and variances, $\boldsymbol{\theta} = \{\mathbf{W}_l, \mathbf{b}_l, \boldsymbol{\Sigma}_l\}_{l=1}^{L+1}$. Though we have modeled the activations as random variables, we can recover standard, feed-forward neural networks when $\boldsymbol{\Sigma}_l \rightarrow \mathbf{0}$ for $l = 1, \dots, L$.

Our goal is to infer the posterior distribution over parameters and activations given a set of observed inputs and outputs,

$$p(\boldsymbol{\theta}, \{\mathbf{a}_n\}_{n=1}^N \mid \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N) \propto p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{y}_n, \mathbf{a}_{n,1}, \dots, \mathbf{a}_{n,L} \mid \mathbf{x}_n, \boldsymbol{\theta}), \quad (3)$$

under a prior $p(\boldsymbol{\theta})$.

MCMC methods like Hamiltonian Monte Carlo and variational approaches like Bayes by back-prop [Blundell et al., 2015] could be used for posterior inference since they only require that the log probability be differentiable. Instead, we introduce additional structure to the model which allows for alternative methods that may yield more efficient inference. Specifically, we propose a deep generative model that relaxes the model above so that it is amenable to Pólya-gamma augmentation [Polson et al., 2013]. This renders the model conditionally linear and Gaussian. Given this, with tractable conditional distributions over large subsets of variables, methods such as Gibbs sampling can update many variables at once and potentially converge more quickly than local, gradient-based algorithms.

We call the model a *bow tie network* due to the shape of the conditional distribution of its activation function. After introducing the model, we derive a Gibbs sampling algorithm that leverages the simple conditionals to perform Bayesian inference. We demonstrate a prototype of the model and

algorithm on various regression tasks and compare it to alternative algorithms for Bayesian deep learning.

2 Bow Tie Networks

Take the generative model from eq. 1 and assume that $f(u) = \max\{0, u\}$ is the rectified linear (ReLU) function. Nonlinear activation functions such as this make posterior inference hard—if f were the identity, we could place conjugate priors on the parameters and derive closed form Gibbs updates for the parameters and activations (more details below). Observe, however, that we could equivalently write the rectified linear function as $f(u) = zu$ where $z = \mathbb{I}[u > 0]$ is a binary activation that determines whether the node is *on* ($z = 1$) or *off* ($z = 0$). Intuitively, if we knew the binary activations, then f would be conditionally linear, which would make inference much simpler.

Motivated by this observation, we propose a stochastic relaxation of the rectified linear network. In the notation of eq. 1, let

$$\mu_{n,l} \triangleq \mathbf{z}_{n,l} \odot \mathbf{u}_{n,l} \quad (4)$$

$$\mathbf{z}_{n,l} \stackrel{\text{ind}}{\sim} \text{Bern}(\sigma(\mathbf{u}_{n,l}/\tau)), \quad (5)$$

where \odot denotes the elementwise product, σ is the logistic function, $\mathbf{w}_{l,d} \in \mathbb{R}^{D_{l-1}}$ is the d -th row of \mathbf{W}_l , and $\tau \geq 0$ is a temperature parameter. Under this model, the rectified linear units are turned on or off stochastically depending on their input.

Figure 1 shows the conditional distribution of the activation $a_{n,l,d}$ given the input $u_{n,l,d}$ and marginalizing over the binary activations $z_{n,l,d}$, for a few values of the temperature τ and the variance $\eta_{l,d}^2 = [\Sigma_l]_{dd}$. For intermediate values of the temperature and noise, the conditional distribution looks like a bow tie, hence the name.

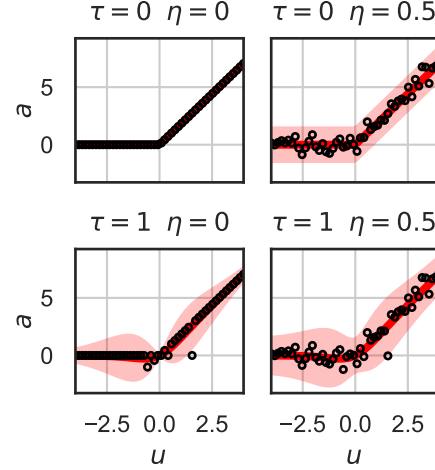


Figure 1: Conditional distribution of the activation a given the input u (marginalizing over the binary activations) for various settings of the temperature τ and the noise η . Red: the conditional mean and variance; black dots: samples from the distribution. When $\tau = \eta = 0$, we recover the standard rectified linear function; for nonzero values, the conditional distribution looks like a bow tie.

3 Bayesian inference via Pólya-gamma augmentation

The bow tie network introduces binary activations to the model so that,

$$p(\mathbf{y}_n, \mathbf{a}_n, \mathbf{z}_n \mid \mathbf{x}_n, \boldsymbol{\theta}) = \left[\prod_{l=1}^L \text{Bern}(\mathbf{z}_{n,l} \mid \sigma(\mathbf{u}_{n,l}/\tau)) \mathcal{N}(\mathbf{a}_{n,l} \mid \mathbf{z}_{n,l} \odot \mathbf{u}_{n,l}, \Sigma_l) \right] \mathcal{N}(\mathbf{y}_n \mid \mathbf{u}_{n,L+1}, \Sigma_{L+1}) \quad (6)$$

The Bernoulli terms still make posterior inference challenging. However, this formulation is amenable to Pólya-gamma (PG) augmentation [Polson et al., 2013]—an auxiliary variable method that equates the Bernoulli pmf with a scale-mean mixture of Gaussians. Let $\boldsymbol{\gamma}_n = \{\gamma_{n,l}\}_{l=1}^L$ with $\gamma_{n,l} \in \mathbb{R}_+^{D_l}$ denote the collection of PG auxiliary variables. The joint distribution on the extended space is,

$$p(\mathbf{y}_n, \mathbf{a}_n, \mathbf{z}_n, \boldsymbol{\gamma}_n \mid \mathbf{x}_n, \boldsymbol{\theta}) \propto \left[\prod_{l=1}^L \text{PG}(\gamma_{n,l} \mid 1, 0) e^{\tau^{-1}(\mathbf{z}_{n,l} - \frac{1}{2})^\top \mathbf{u}_{n,l} - \frac{1}{2\tau^2} \mathbf{u}_{n,l}^\top \mathbf{G}_{n,l} \mathbf{u}_{n,l}} \mathcal{N}(\mathbf{a}_{n,l} \mid \mathbf{z}_{n,l} \odot \mathbf{u}_{n,l}, \Sigma_l) \right] \times \mathcal{N}(\mathbf{y}_n \mid \mathbf{u}_{n,L+1}, \Sigma_{L+1}), \quad (7)$$

where $\text{PG}(\gamma_{n,l} \mid 1, 0)$ is the probability density of a bank of D_l independent, standard PG random variables and $\mathbf{G}_{n,l} = \text{diag}(\gamma_{n,l})$.

Dataset	Test RMSE			Test Log-likelihood		
	BBB	fBNN	Bow Tie	BBB	fBNN	Bow Tie
Boston	3.171 ± 0.149	2.378 ± 0.104	3.324 ± 0.167	-2.602 ± 0.031	-2.301 ± 0.038	-2.506 ± 0.035
Concrete	5.678 ± 0.087	4.935 ± 0.180	5.241 ± 0.206	-3.149 ± 0.018	-3.096 ± 0.016	-3.045 ± 0.041
Energy	0.565 ± 0.018	0.412 ± 0.017	1.096 ± 0.086	-1.500 ± 0.006	-0.684 ± 0.020	-1.334 ± 0.060
Wine	0.643 ± 0.012	0.673 ± 0.014	0.651 ± 0.020	-0.977 ± 0.017	-1.040 ± 0.013	-0.992 ± 0.059
Yacht	1.174 ± 0.086	0.607 ± 0.068	0.896 ± 0.083	-2.408 ± 0.007	-1.033 ± 0.033	-1.139 ± 0.093

Table 1: Performance of bow tie networks, Bayes by backprop (BBB) Blundell et al. [2015], and functional variational Bayesian neural networks (fBNN) Sun et al. [2019] on UCI regression datasets Dua and Graff [2017]. Bow tie networks are generally competitive with both of these methods. BBB and fBNN results are from Sun et al. [2019].

The log probability under the augmented model is quadratic in $\mathbf{u}_{n,l}$, which is in turn a linear function of the weights, biases, and preceding layer’s activations. Thus, the augmented model is conjugate with a Gaussian prior on the weights and an inverse Wishart prior on the covariance matrices (or more generally a matrix-normal inverse Wishart prior). Moreover, the conditional distribution of the activations $\mathbf{a}_n = (\mathbf{a}_{n,1}, \dots, \mathbf{a}_{n,L})$ is a linear Gaussian chain. That means we can Gibbs sample the activations via an efficient forward filtering backward sampling algorithm in order $O(NLD^3)$ time (where $D = \max_l D_l$). When the covariances are constrained to be diagonal, the binary activations follow simple, independent Bernoulli conditionals. Likewise, the PG variables follow a tilted PG conditional distribution, which is also easy to sample from. Finally, note that the latent variables are conditionally independent across data points, so the Gibbs sampler can be easily parallelized across this dimension. Complete details are in the appendix.

4 Experiments

We evaluated bow tie networks on a standard set of regression tasks from the UCI dataset repository Dua and Graff [2017]. Following Sun et al. [2019], we selected only datasets with fewer than 2000 data points so that full-batch learning was tractable. The bow tie network architecture used a single hidden layer of 50 units, consistent with the experiment in Sun et al. [2019]. The activation temperature τ was fixed at 0.1 and 1,000 samples were collected after a burn-in of 26,000 steps. See table 1 for the experimental results. Bow tie networks were generally competitive with both Bayes by backprop (BBB) and functional variational Bayesian neural networks (fBNNs) on the regression tasks.

5 Conclusions

We presented bow tie networks, a novel model for Bayesian neural networks that allows for fast and accurate inference via Pólya-gamma augmentation. In the future we will improve the performance of bow tie networks by exploring alternative sampling strategies such as annealing the temperature over training. We also hope to scale bow tie networks to larger datasets and explore connections with inference in the function space.

References

- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1613–1622, Lille, France, 07–09 Jul 2015. PMLR.
- Nicholas G Polson, James G Scott, and Jesse Windle. Bayesian inference for logistic models using Pólya-gamma latent variables. *Journal of the American statistical Association*, 108(504): 1339–1349, 2013.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Shengyang Sun, Guodong Zhang, Jiaxin Shi, and Roger Grosse. Functional variational bayesian neural networks. *arXiv preprint arXiv:1903.05779*, 2019.

A Complete conditional distributions

In this section we provide the complete conditional distributions necessary for the Gibbs sampling algorithm described in Section 3.

Let $\tilde{\mathbf{a}}_{n,l} = (\mathbf{a}_{n,l}, 1)$ and $\tilde{\mathbf{w}}_{l,d} \triangleq (\mathbf{w}_{l,d}, b_{l,d})$ where $\mathbf{w}_{l,d} \in \mathbb{R}^{D_l-1}$ is the d -th row of the weights \mathbf{W}_l . Assume, for simplicity, that the noise covariances are diagonal $\Sigma_l = \text{diag}(\eta_{l,1}^2, \dots, \eta_{l,D_l}^2)$. Then assume a factored prior,

$$p(\boldsymbol{\theta}) = \prod_{l=1}^{D_l} \prod_{d=1} \mathcal{N}(\tilde{\mathbf{w}}_{l,d} \mid \boldsymbol{\mu}_0, \Sigma_0) \text{Ga}(\eta_{l,d}^{-2} \mid \alpha_0, \beta_0). \quad (8)$$

Each of the updates below can be performed in parallel. For example, activations can be sampled in parallel for all data points, the weights can be sampled in parallel for indices all layers and output dimensions, and the binary activations and auxiliary variables can be sampled in parallel for all data points, layers, and dimensions. Leveraging this parallelism could lead to much faster implementations on massively parallel architectures like GPUs.

Conditional distribution of the activations. The most interesting conditional distribution is that of the activations. Given the parameters, binary activations, and auxiliary variables, the conditional distribution of the activations is proportional to,

$$p(\mathbf{a}_n \mid \mathbf{y}_n, \mathbf{z}_n, \boldsymbol{\gamma}_n, \mathbf{x}_n, \boldsymbol{\theta}) \propto \exp \left\{ -\frac{1}{2} \sum_{l=1}^L \mathbf{a}_{n,l}^\top \mathbf{J}_{n,l} \mathbf{a}_{n,l} - \sum_{l=2}^L \mathbf{a}_{n,l}^\top \mathbf{L}_{n,l} \mathbf{a}_{n,l-1} + \sum_{l=1}^L \mathbf{h}_{n,l}^\top \mathbf{a}_{n,l} \right\}. \quad (9)$$

This is a linear Gaussian chain, just like a Gaussian linear dynamical system (LDS). To simplify the expressions below, let $\mathbf{Z}_{n,l} = \text{diag}(\mathbf{z}_{n,l})$. Then the coefficients are,

$$\mathbf{J}_{n,l} = \Sigma_l^{-1} + \mathbf{W}_{l+1}^\top (\tau^{-2} \mathbf{G}_{n,l+1} + \mathbf{Z}_{n,l+1} \Sigma_{l+1}^{-1} \mathbf{Z}_{n,l+1}) \mathbf{W}_{l+1} \quad \text{for } l < L \quad (10)$$

$$\mathbf{J}_{n,L} = \Sigma_L^{-1} + \mathbf{W}_{L+1}^\top \Sigma_{L+1} \mathbf{W}_{L+1} \quad (11)$$

$$\mathbf{L}_{n,l} = -\mathbf{Z}_{n,l} \Sigma_l^{-1} \mathbf{W}_l \quad (12)$$

$$\mathbf{h}_{n,1} = \mathbf{Z}_1 \Sigma_1^{-1} \mathbf{b}_1 + \mathbf{W}_2^\top \left(\frac{1}{\tau} (\mathbf{z}_{n,2} - \frac{1}{2}) - \frac{1}{\tau^2} \mathbf{G}_{n,2} \mathbf{b}_2 - \mathbf{Z}_{n,2} \Sigma_2^{-1} \mathbf{b}_2 \right) + \mathbf{Z}_{n,1} \Sigma_1^{-1} \mathbf{W}_1 \mathbf{x}_n \quad (13)$$

$$\mathbf{h}_{n,l} = \mathbf{Z}_l \Sigma_l^{-1} \mathbf{b}_l + \mathbf{W}_{l+1}^\top \left(\frac{1}{\tau} (\mathbf{z}_{n,l+1} - \frac{1}{2}) - \frac{1}{\tau^2} \mathbf{G}_{n,l+1} \mathbf{b}_{l+1} - \mathbf{Z}_{n,l+1} \Sigma_{l+1}^{-1} \mathbf{b}_{l+1} \right) \quad (14)$$

$$\mathbf{h}_{n,L} = \mathbf{Z}_L \Sigma_L^{-1} \mathbf{b}_L + \mathbf{W}_{L+1}^\top \Sigma_{L+1}^{-1} (\mathbf{y}_n - \mathbf{b}_{L+1}) \quad (15)$$

We can draw exact samples from this conditional distribution in $O(LD^3)$ time where $D = \max_l D_l$ using the forward filtering backward sampling algorithm. Note that we have presented the conditional distribution in “information form” rather than in terms of mean parameters, as linear Gaussian dynamical systems.

Conditional distribution of the weights and biases. Under this prior, the conditional distribution of the weights is,

$$p(\tilde{\mathbf{w}}_{l,d} \mid \mathbf{a}_{n,l}, \mathbf{z}_{n,l}, \boldsymbol{\gamma}_{n,l}, \mathbf{a}_{n,l-1}) = \mathcal{N}(\tilde{\mathbf{w}}_{l,d} \mid \mathbf{J}_{l,d}^{-1} \mathbf{h}_{l,d}, \mathbf{J}_{l,d}^{-1}) \quad (16)$$

$$\text{where } \mathbf{J}_{l,d} = \Sigma_0^{-1} + \sum_{n=1}^N (\tau^{-2} \gamma_{n,l,d} + \eta_{l,d}^{-2} z_{n,l,d}) \tilde{\mathbf{a}}_{n,l} \tilde{\mathbf{a}}_{n,l}^\top \quad (17)$$

$$\mathbf{h}_{l,d} = \Sigma_0^{-1} \boldsymbol{\mu}_0 + \sum_{n=1}^N \left(\tau^{-1} (z_{n,l,d} - \frac{1}{2}) + \eta_{l,d}^{-2} z_{n,l,d} \mathbf{a}_{n,l,d} \right) \tilde{\mathbf{a}}_{n,l} \quad (18)$$

Conditional distribution of the variances. The conditional distribution of the precisions (inverse variances) is,

$$p(\eta_{l,d}^{-2} \mid \mathbf{a}_{n,l}, \mathbf{z}_{n,l}, \boldsymbol{\gamma}_{n,l}, \mathbf{a}_{n,l-1}, \mathbf{w}_{l,d}, b_{l,d}) = \text{Ga}(\eta_{l,d}^{-2} \mid \alpha_{l,d}, \beta_{l,d}) \quad (19)$$

$$\text{where } \alpha_{l,d} = \alpha_0 + \frac{N}{2} \quad (20)$$

$$\beta_{l,d} = \beta_0 + \frac{1}{2} \sum_{n=1}^N (a_{n,l,d} - z_{n,l,d} u_{n,l,d})^2 \quad (21)$$

Conditional distribution of the binary activations. The conditional distribution of the binary activations, marginalizing over the PG auxiliary variables, is,

$$p(z_{n,l,d} \mid \mathbf{a}_{n,l}, \mathbf{a}_{n,l-1}, \boldsymbol{\theta}) = \text{Bern}(\sigma(\nu_{n,l,d})) \quad (22)$$

$$\text{where } \nu_{n,l,d} = \frac{u_{n,l,d}}{\tau} + \frac{a_{n,l,d}u_{n,l,d}}{\eta_{l,d}^2} - \frac{u_{n,l,d}^2}{2\eta_{l,d}^2} \quad (23)$$

Conditional distribution of the PG auxiliary variables. Since we Gibbs sample the binary activations from their conditional distribution under the non-augmented model (eq. 6), we must immediately Gibbs sample the PG auxiliary variables after updating the binary activations. Their conditional distribution is,

$$p(\gamma_{n,l,d} \mid \tilde{\mathbf{a}}_{n,l-1}, \boldsymbol{\theta}) = \text{PG}(\gamma_{n,l,d} \mid 1, \tau^{-1}u_{n,l,d}) \quad (24)$$