
Data Augmentation for Bayesian Deep Learning

Yuexi Wang,

Joint work with *Nick Polson and Vadim Sokolov*

Booth School of Business
University of Chicago

June 27th, ISBA 2022

Bayesian Deep Learning

Empirical success of deep learning amply illustrated.

Theoretical developments

- approximability
 - Why and when neural networks generalize well?
 - When do deep networks out-perform shallow ones?
 - Which activation functions and with how many layers?
- rate of convergence, including both frequentist and Bayesian point of view

However, training deep learners is still challenging

- high dimensional search space
- non-convex objective function

Issues like local traps, miscalibration and overfitting.

We provide a *Bayesian alternative* that could help.

What is Deep Learning (DL)?

- DL reconstructs high-dimensional **input-output mappings**.
- Given training data, *inputs* $\mathbf{x}_i \in \mathbb{R}^p$ and *outputs* $y_i = f(\mathbf{x}_i)$ for $1 \leq i \leq n$, **learn** $f_{\hat{\theta}}$ such that

$$f_{\hat{\theta}}(\mathbf{x}) \approx f(\mathbf{x}) \quad \text{for} \quad \mathbf{x} \notin \{\mathbf{x}_i\}_{i=1}^n.$$

- Training neural networks is then positioned as an *optimization problem* for finding values $\hat{\theta} \in \mathbb{R}^T$ that minimize regularized empirical risk

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n \ell(y_i, f_{\theta}(\mathbf{x}_i)) + \phi(\theta) \quad (1)$$

where $\phi(\theta)$ is a penalty over the weights and offset parameters.

In practice, this is most often carried out with some form of **stochastic gradient descent (SGD)** (see e.g. Polson and Sokolov (2017) for an overview).

Deep Learning Mappings

A deep neural network $f_{\theta}(\mathbf{x})$ is an iterative mapping specified by *hierarchical layers of abstraction*.

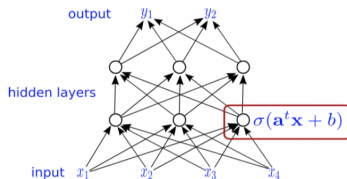


Figure 1: Representation as a direct graph of a network with two hidden layers $L = 2$ and width vector $\mathbf{p} = (4, 3, 3, 2)$.

With $L \in \mathbb{N}$ we denote the *number of hidden layers* and with $p_l \in \mathbb{N}$ the number of neurons at the l^{th} layer.

Setting $p_0 = p$ and $p_{L+1} = 1$, we denote with $\mathbf{p} = (p_0, \dots, p_{L+1})' \in \mathbb{N}^{L+2}$ the vector of neuron counts for the entire network.

We denote the parameters with $\theta = \{(\mathbf{W}_0, b_0), (\mathbf{W}_1, b_1), \dots, (\mathbf{W}_L, b_L)\}$

Regularization DL

*It is commonly agreed that generalizability of neural networks can be improved with **regularization**.* (Goodfellow et al., 2016)

- Regularization through ReLU activation.
- **Dropout regularization** (Srivastava et al., 2014) samples from (and averages over) thinned networks obtained by *randomly* dropping out nodes together with their connections.
- While motivated as stochastic regularization, dropout can be regarded as ℓ_2 *regularization* obtained by margining out dropout noise (Wager, 2014).
- Conceptually related to **Bayesian model averaging** (ℓ_0 penalization).

How does Bayesian framework come into play?

Duality between Bayesian Simulation and Regularization

We can specify different amount of regularization $\{\lambda_l, \phi_l(\cdot)\}$ for each layer.

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) + \sum_{l=0}^L \lambda_l \phi_l(W_l, b_l). \quad (2)$$

Solving the problem is equivalent to find the *maximum a posteriori* (MAP) estimate to the following posterior

$$p(\theta | \mathbf{y}) = p(\mathbf{y} | \theta)p(\theta)/p(\mathbf{y}),$$

where

$$p(\mathbf{y}|\theta) \propto \exp \left\{ - \sum_{i=1}^n \ell(y_i, f_{\theta}(\mathbf{x}_i)) \right\}, \quad p(\theta) \propto \exp \left\{ - \sum_{l=0}^L \lambda_l \phi_l(W_l, b_l) \right\}.$$

Here $p(\theta)$ can be interpreted as a prior probability distribution and the log-prior as the regularization penalty.

From Optimization to Sampling

By exploiting the duality, we can convert the composite functions

$$\begin{aligned}\mathbf{y} &= f_0(Z_1 W_0 + b_0), \\ Z_l &= f_l(Z_{l+1} W_l + b_l), \quad l = 1, \dots, L, \\ Z_{L+1} &= \mathbf{x},\end{aligned}\tag{3}$$

to their stochastic resemble as

$$\begin{aligned}\mathbf{y} \mid Z_1 &\sim p(\mathbf{y} \mid Z_1), \\ Z_l &\sim N(f_l(W_l Z_{l+1} + b_l), \tau_l^2), \quad l = 1, 2, \dots, L, \\ Z_{L+1} &= \mathbf{x}.\end{aligned}$$

Now the hidden variables $Z = (Z_1, \dots, Z_L)'$ can be viewed as data augmentation variables and be updated via sampling.

A Stochastic Top Layer

For the ease of computation, we only replace the top layer of DNNs with a stochastic layer.

We implement the solutions with iterative search.

- 1 **DA-update** for the top layer W_0, b_0

$$p(W_0, b_0 | Z_1^{(t)}, \mathbf{y}) \propto \exp \left\{ -\frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i) | \mathbf{B}^{(t)}) + \lambda_0 \phi_0(W_0, b_0) \right\}$$

- 2 **SGD-update** for the deep architecture \mathbf{B}

$$\mathbf{B}^{(t+1)} = \arg \min_{\mathbf{B}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i) | (W_0, b_0)^{(t+1)}) + \sum_{l=1}^L \lambda_l \phi_l(W_l, b_l)$$

- 3 Sample $Z_1^{(t+1)}$ from a normal distribution $\mathcal{N}(\mu_z^{(t)}, \sigma_z^{(t)})$ (determined jointly by $\{\theta^{(t)}, \mathbf{x}, \mathbf{y}\}$).

DA for Common Activation Functions

Data augmentation tricks allow us to express the likelihood as an expectation of a **weighted L^2 -norm** or **mixture of normals**

$$\exp \left\{ -\ell(\mathbf{y}, f_{\theta}(\mathbf{x})) \right\} = \int_0^{\infty} \exp \left(-Q(\mathbf{y} \mid f_{\theta}(\mathbf{x}), \omega) \right) p(\omega) d\omega,$$

where $p(\omega)$ is the prior on the auxiliary variables ω and $Q(\cdot)$ is designed to be a quadratic form.

Commonly used activation functions such as ReLU, logit, lasso and check can be expressed in the form (Polson et al., 2013)

$$\exp(-\max(1-x, 0)) = E_{\omega} \left\{ \frac{1}{\sqrt{2\pi\omega}} \exp \left(-\frac{1}{2\omega} (x-1-\omega)^2 \right) \right\}, \quad \text{where } \omega \sim \mathcal{GIG}(1, 0, 0),$$

$$\exp(-\log(1 + e^x)) = E_{\omega} \left\{ \exp \left(-\frac{1}{2} \omega x^2 \right) \right\}, \quad \text{where } \omega \sim \mathcal{PG}(1, 0),$$

$$\exp(-|x|) = E_{\omega} \left\{ \frac{1}{\sqrt{2\pi\omega}} \exp \left(-\frac{1}{2\omega} x^2 \right) \right\}, \quad \text{where } \omega \sim \mathcal{E}\left(\frac{1}{2}\right).$$

\mathcal{GIG} – the Generalized Inverse Gaussian distribution

\mathcal{PG} – Pólya Gamma

\mathcal{E} – exponential distribution.

Solutions to Data Augmentation

Two approaches to handle DA (Geyer, 1996)

- missing data methods like Expectation-Maximization (EM)
 - 😊 based on a probabilistic *approximation* of the objective function
 - 😞 less concerned with exploring Θ
- stochastic search methods like MCMC
 - 😊 visiting the entire range of Θ
 - 😞 less tied to the properties of the function

EM Algorithms

Consider constructing a surrogate optimization problem (Lange et al., 2000) which has the same solution to original problem

$$H(\theta) = \mathbb{E}_{\omega|\theta} \left[\exp \left(-\frac{1}{n} \sum_{i=1}^n Q(y_i | f_{\theta}(\mathbf{x}_i), \omega) - \sum_{l=0}^L \lambda_l \phi_l(\mathbf{w}_l, \mathbf{b}_l) \right) \right],$$

which is a concave function to be maximized.

Using Jensen's inequality we construct

$$G(\theta | \theta^{(t)}) = -\mathbb{E}_{\omega|\theta^{(t)}} \left[\frac{1}{n} \sum_{i=1}^n Q(y_i | f_{\theta}(\mathbf{x}_i), \omega) + \sum_{l=0}^L \lambda_l \phi_l(\mathbf{w}_l, \mathbf{b}_l) \right],$$

and the minorization is satisfied as

$$\log H(\theta) \geq G(\theta | \theta^{(t)}).$$

Maximizing $G(\theta | \theta^{(t)})$ with respect to θ drives $H(\theta)$ uphill.

EM in Action: Logistic Regression

Adopting the logistic regression model from Polson and Scott (2013), we focus on the penalization of W_0 as

$$\hat{W}_0 = \arg \min_{W_0} \left[\frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp \left(- y_i f_{\mathbf{B}}(\mathbf{x}_i) W_0 \right) \right) + \phi(W_0 | \tau) \right],$$

The outcomes y_i are coded as ± 1 , and τ is assumed fixed.

For likelihood function ℓ and regularization penalty ϕ , we assume

$$p(y_i | \sigma) \propto \int_0^\infty \frac{\sqrt{\omega_i}}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{\omega_i}{2\sigma^2} \left(y_i f_{\mathbf{B}}(\mathbf{x}_i) W_0 - \frac{1}{2\omega_i} \right)^2 \right\} p(\omega_i) d\omega_i,$$
$$p(W_0 | \tau) = \int_0^\infty \frac{\sqrt{\lambda}}{\sqrt{2\pi\tau}} \exp \left\{ -\frac{\lambda}{2\tau^2} (W_0 - \mu_W - \kappa_W \lambda^{-1})^2 \right\} p(\lambda) d\lambda,$$

where λ follows a Pólya distribution prior $P(\lambda)$.

The deterministic updates of $\{W_0, \omega, \lambda\}$ generate a sequence of estimates that converges to a stationary point of posterior

MCMC and J-copies

A sequence can be simulated using Gibbs conditionals,

$$\begin{aligned}p(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\omega}^{(t)}, \mathbf{y}) &\propto \exp\left(-Q(\mathbf{y} \mid f_{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\omega}^{(t)})\right)p(\boldsymbol{\theta}), \\p(\boldsymbol{\omega}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}, \mathbf{y}) &\propto \exp\left(-Q(\mathbf{y} \mid f_{\boldsymbol{\theta}^{(t)}}(\mathbf{x}), \boldsymbol{\omega})\right)p(\boldsymbol{\omega}).\end{aligned}$$

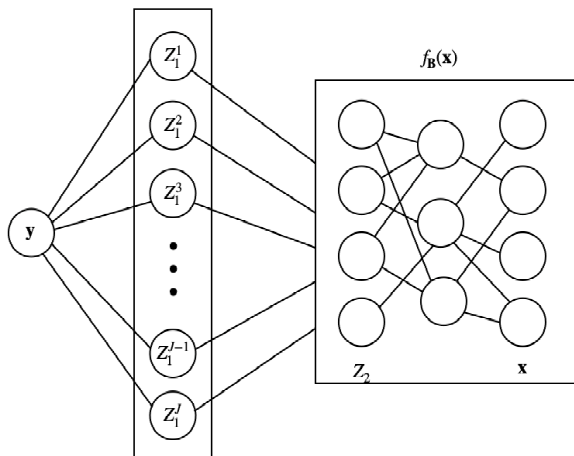
Inspired by the simulated annealing algorithm (Metropolis et al., 1953), we add a scaling factor J to help avoiding the trapping attraction of local maxima

$$\pi_J(\boldsymbol{\theta}) = \exp\{-Jf(\boldsymbol{\theta})\} / Z_J \text{ for } \boldsymbol{\theta} \in \Theta \quad (4)$$

where $Z_J = \int_{\Theta} \exp\{-Jf(\boldsymbol{\theta})\} d\boldsymbol{\theta}$ is an appropriate normalizing constant.

MCMC and J-copies

To simulate the posterior mode without evaluating the likelihood directly (Jacquier et al., 2007), we sample J independent copies of hidden variable Z_1 .



$$Z_1^j \sim Z_1 | y, \mathbf{x}, \mathbf{B}, \mathbf{W}_0$$

MCMC in Action: Gaussian Regression

We consider the regression model as

$$y_i = z_{1,i} W_0 + b_0 + \epsilon_{0,i}, \quad \text{where } y_i \in (-\infty, \infty), \epsilon_{0,i} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \tau_0^2),$$
$$z_{1,i} = f_{\mathbf{B}}(\mathbf{x}_i) + \epsilon_{z,i}, \quad \text{where } \epsilon_{z,i} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \tau_z^2).$$

The posterior updates are given by

$$\hat{W}_0 = \text{Cov}(Z_1, \mathbf{y}) / \text{Var}(Z_1),$$

$$\hat{b}_0 = \bar{\mathbf{y}} - W_0 \bar{Z}_1$$

$$p(Z_1 \mid \mathbf{y}, \mathbf{x}, \theta) = C_z \exp \left\{ -\frac{1}{2\tau_0^2} \|\mathbf{y} - Z_1 W_0 - b_0\|^2 - \frac{1}{2\tau_z^2} \|Z_1 - f_{\mathbf{B}}(\mathbf{x})\|^2 \right\},$$

where $\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n y_i$ and C_z is a normalizing constant.

MCMC in Action: SVM

Polson and Scott (2011) and Mallick et al. (2005) write the support vector machine model as

$$\mathbf{y} = \mathbf{Z}_1 \mathbf{W}_0 + \boldsymbol{\lambda} + \sqrt{\boldsymbol{\lambda}} \boldsymbol{\epsilon}_0, \text{ where } \boldsymbol{\lambda} \sim p(\boldsymbol{\lambda}),$$

where $p(\boldsymbol{\lambda})$ follows a flat uniform prior.

By incorporating the augmentation variable $\boldsymbol{\lambda}$, the ReLU deep learning model can be written as

$$y_i = z_{1,i} \mathbf{W}_0 + \lambda_i + \sqrt{\lambda_i} \epsilon_{0,i}, \text{ where } y_i \in \{-1, 1\}, \epsilon_{0,i} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \tau_0^2), \\ z_{1,i} = f_{\mathbf{B}}(\mathbf{x}_i) + \epsilon_{z,i}, \text{ where } \epsilon_{z,i} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \tau_z^2).$$

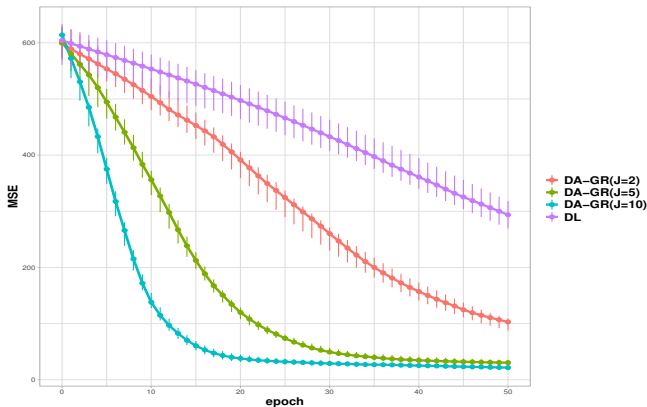
In order to generate the latent variables, we use conditional Gibbs sampling as

$$\begin{aligned} \lambda_i^{-1} \mid \mathbf{W}_0, y_i, z_{1,i} &\sim \mathcal{IG}(|1 - y_i z_{1,i} \mathbf{W}_0|^{-1}, \tau_0^{-2}) \\ \mathbf{W}_0 \mid \mathbf{y}, \mathbf{Z}_1, \boldsymbol{\lambda} &\sim \mathcal{N}(\mu_w, \sigma_w^2) \\ \mathbf{Z}_1 \mid \mathbf{y}, \mathbf{x}, \mathbf{W}_0, \mathbf{B} &\sim \mathcal{N}(\mu_z, \sigma_z^2) \end{aligned}$$

Empirical Analysis

Boston Housing Dataset

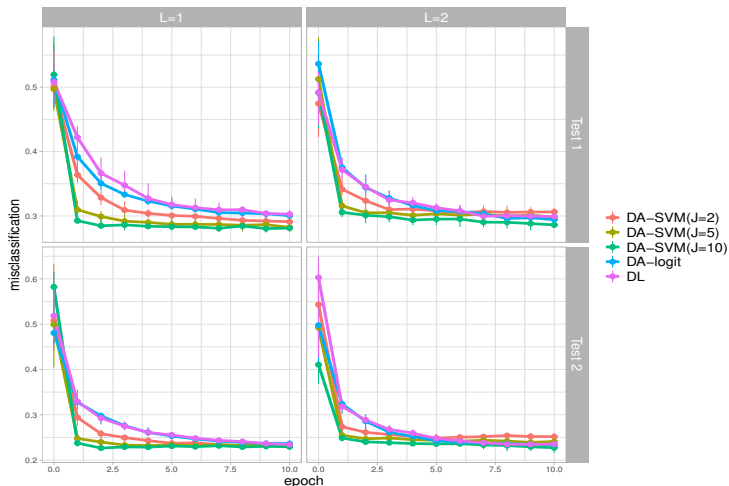
Boston Housing dataset¹ has $n = 506$ observations with 13 features. We adopt a one-layer ReLU networks with 64 units and set the dropout rate to be 0.5.



¹<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>

Wine Quality Dataset

The Wine Quality Data Set ² contains 4 898 observations with 11 features. Two types of binary classifications are considered.



²P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis, 'Wine Quality Data Set', UCI Machine Learning Repository.

Concluding Remarks

- We proposed a new approach to update the neural network hidden nodes via sampling.
- The objective function can be rewritten in **weighted least squares** form using data augmentation, which can be implemented at scale with accelerated linear algebra methods.
- We provide motivating examples for commonly used activation functions. Substantial performance gains are observed in empirical analysis.

Thank you!

Reference: Wang, Y., Polson, N. G., & Sokolov, V. O. (2019). Scalable Data Augmentation for Deep Learning.