

---

# Variational Inference and Dropout

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1

## 2 1 Variational Inference for Latent Variable Model

3 Let  $x = (x_i, i \in \llbracket 1, n \rrbracket)$  and  $y = (y_i, i \in \llbracket 1, n \rrbracket)$  be i.i.d. input-output pairs and  $w \in W \subseteq \mathbb{R}^J$  be a  
4 latent variable. The joint distribution of  $y, w$  can be written as:

$$p(y, w|x) = p(w) \prod_{i=1}^n p_i(y_i|x_i, w) \quad (1)$$

5 Our goal is to compute the posterior distribution  $p(w|y, x)$  given the input-output pairs  $x, y$ . The  
6 variational Inference (VI) algorithm [?] consists of minimizing the KL divergence between a can-  
7 didate family of parametric distributions  $\{q(w; \theta), \theta \in \Theta\}$  and the posterior distribution  $p(w|y, x)$   
8 of the global latent variable  $w$ . For example,  $q(w; \theta)$  belongs to a simple family of distributions  
9 such as the multivariate Gaussian family with mean  $\rho$  and covariance matrix  $\sigma^2 \mathbf{I}$ , where we have  
10  $\theta = (\rho, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+^*$ . VI can be framed as an optimization problem, usually in terms of KL  
11 divergence, of the following form:

$$\theta^* = \arg \min_{\theta \in \Theta} \text{KL}(q(w; \theta) || p(w|y, x)) = \arg \min_{\theta \in \Theta} \mathcal{L}(\theta) \quad (2)$$

12 where  $\mathcal{L}(\theta) := n^{-1} \sum_{i=1}^n \mathcal{L}_i(\theta)$  for  $\theta \in \Theta$  with :

$$\mathcal{L}_i(\theta) := - \int_W q(w; \theta) \log p_i(y_i|x_i, w) dw + \frac{1}{n} \text{KL}(q(w; \theta) || p(w)) = r_i(\theta) + d(\theta), \quad (3)$$

13 Directly optimizing the finite sum objective (2) can be infeasible especially when  $n \gg 1$ . This is  
14 because evaluating the objective function  $\mathcal{L}(\theta)$  requires a full pass of computation over the entire  
15 dataset, and this approach does not adapt to complex models when the last integral cannot be eval-  
16 uated analytically. For instance, such optimization problem is notoriously hard for Bayesian neural  
17 networks [?].

18 To apply the MISSO method with a stochastic surrogate model (??), we consider the following  
19 quadratic surrogate function. For any  $i \in \llbracket 1, n \rrbracket$ , we take:

$$\hat{\mathcal{L}}_i(\theta; \bar{\theta}) := \mathcal{L}_i(\bar{\theta}) + \langle \nabla \mathcal{L}_i(\bar{\theta}) | \theta - \bar{\theta} \rangle + \frac{L}{2} \|\bar{\theta} - \theta\|^2 \quad (4)$$

20 where  $L$  is the smoothness modulus of  $\mathcal{L}_i$  at  $\bar{\theta}$ . To compute the gradient  $\nabla \mathcal{L}_i(\bar{\theta})$ , we apply the  
21 re-parametrization technique suggested in [???]. Let  $t : \Theta \times \mathbb{R}^d \mapsto \mathbb{R}^d$  be a measurable function  
22 and  $\phi$  be the density of the standard normal distribution  $\mathcal{N}_d(0, \mathbf{I})$ . The function  $t$  is designed such  
23 that for all  $\bar{\theta} \in \Theta$  and for  $\epsilon \sim \phi(\cdot)$ , the distribution of the random vector  $W = t(\bar{\theta}, \epsilon)$  is the same as  
24  $q(\cdot, \bar{\theta})$ . It follows from [?, Proposition 1] that:

$$\nabla \int_W \log p_i(y_i|x_i, w) q(w, \bar{\theta}) dw = \int_{\mathbb{R}^d} \mathbf{J}_t^\theta(\bar{\theta}; e) \nabla \log p_i(y_i|x_i, t(\bar{\theta}, e)) \phi(e) de \quad (5)$$

where for each  $e \in \mathbb{R}^d$ ,  $J_t^\theta(\bar{\theta}; e)$  is the Jacobian of the function  $t(\cdot, e)$  with respect to  $\theta$  evaluated at  $\bar{\theta}$ . Consequently, we can apply the MISSO method to tackle the VI problem with the pair  $(r_i(\theta; \bar{\theta}, e), \phi(e))$  where  $r_i(\theta; \bar{\theta}, e)$  reads:

$$r_i(\theta; \bar{\theta}, e) := (-\log p_i(y_i|x_i, t(\bar{\theta}, e)) + d(\bar{\theta})) + \left( -J_t^\theta(\bar{\theta}; e) \nabla \log p_i(y_i|x_i, t(\bar{\theta}, e)) + \nabla d(\bar{\theta}) \right)^\top (\theta - \bar{\theta}) + \frac{L}{2} \|\theta - \bar{\theta}\|^2 \quad (6)$$

## 2 Numerical Experiments

### 2.1 Fitting Bayesian LeNet-5 on MNIST:

In this experiment, we implement the MISSO algorithm for variational inference in the Bayesian variant of LeNet-5 [?] (architecture is described in Appendix .1.1). We train this network on the MNIST dataset [?] used extensively as a benchmark example. The training set is composed of  $N = 55\,000$  handwritten digits,  $28 \times 28$  images. Each image is labelled with its corresponding number (from zero to nine).

Given weight matrices  $(w_\ell)_{\ell=1}^L \in \mathbb{W}^L$ , we put, for each layer  $\ell$  a standard Gaussian prior distributions over those weights:  $p(w_\ell) = \mathcal{N}(0, I)$ , as in [?]. We denote by  $f(x_i, (w_\ell)_{\ell=1}^L)$  the output of the nested function describing the neural network with weights  $(w_\ell)_{\ell=1}^L \in \mathbb{W}^L$  taking as input a data sample  $x_i$ . We assume a softmax likelihood in that classification task:  $p(y_i|x_i, w) = \text{Softmax}(f(x_i, (w_\ell)_{\ell=1}^L))$ .

The variational distribution  $q(w, \theta)$  belongs to the Gaussian multivariate distribution family. Here, the MISSO algorithm coincides with a mini-batch version of the Variational Inference algorithm. At iteration  $k$ , minimizing the sum of stochastic surrogates defined as in (??) and by the quantities (6) yields the following MISSO update: pick a function index  $I_k$  uniformly on  $\llbracket 1, n \rrbracket$ , sample a Monte Carlo batch  $\{e_m^{(k)}\}_{m=1}^{M(k)}$  from the standard Gaussian distribution and update the parameters as  $\theta^{(k)} = \frac{1}{n} \sum_{i=1}^n \theta^{(\tau_i^k)} - \frac{1}{2\gamma} \sum_{i=1}^n \hat{m}_i^k$  where  $\hat{m}_i^k$  are defined recursively as follows:

$$\hat{m}_i^{(k)} \triangleq \begin{cases} -\frac{1}{M(k)} \sum_{m=0}^{M(k)-1} J_\theta^t(e_m^{(k)}) \nabla_\theta \log p_i(y_i, x_i | t(\theta, e_m^{(k)})) + \nabla d(\theta^{(k-1)}) & \text{if } i \in I_k \\ \hat{m}_i^{(k-1)} & \text{otherwise} \end{cases} \quad (7)$$

We compare the convergence behaviors of the following state of the art optimization algorithms, using their vanilla implementations on TensorFlow [?]: the ADAM [?], the Momentum [?] the *Bayes by Back-prop* (BBB) [?] and the Dropout [?] algorithms versus our MISSO update. The loss function (3) and its gradients were computed by Monte Carlo integration using Tensorflow Probability library, based on the reparametrization trick. We use the following hyperparameters for all runs: the learning rate is set to  $10^{-3}$ , we run 100 epochs and use a mini-batch size of 128.

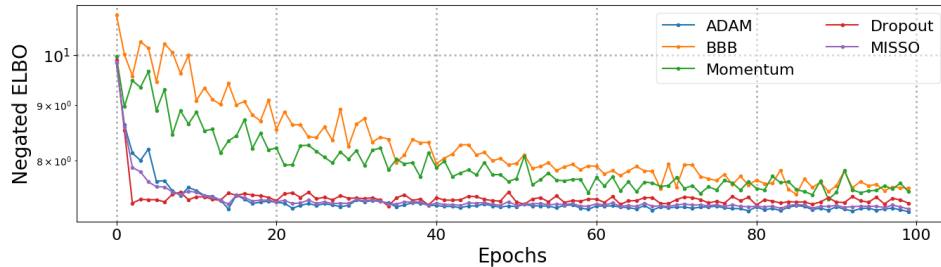


Figure 1: (Incremental Variational Inference) Convergence of the negated ELBO on MNIST.

52 **.1 Incremental Variational Inference for MNIST**

53 **.1.1 Bayesian LeNet-5 Architecture**

layer type	width	stride	padding	input shape	nonlinearity
convolution ( $5 \times 5$ )	6	1	0	$1 \times 32 \times 32$	ReLU
max-pooling ( $2 \times 2$ )		2	0	$6 \times 28 \times 28$	
convolution ( $5 \times 5$ )	6	1	0	$1 \times 14 \times 14$	ReLU
max-pooling ( $2 \times 2$ )		2	0	$16 \times 10 \times 10$	
fully-connected	120			400	ReLU
fully-connected	84			120	ReLU
fully-connected	10			84	

Table 1: LeNet-5 architecture