

Semi-Implicit Variational Inference: Supplementary Material

Mingzhang Yin and Mingyuan Zhou

Algorithm 1 Semi-Implicit Variational Inference (SIVI)

input : Data $\{x_i\}_{1:N}$, joint likelihood $p(\mathbf{x}, \mathbf{z})$, explicit variational distribution $q_\xi(\mathbf{z} | \psi)$ with reparameterization $\mathbf{z} = f(\epsilon, \xi, \psi)$, $\epsilon \sim p(\epsilon)$, implicit layer neural network $T_\phi(\epsilon)$ and source of randomness $q(\epsilon)$

output : Variational parameter ξ for the conditional distribution $q_\xi(\mathbf{z} | \psi)$, variational parameter ϕ for the mixing distribution $q_\phi(\psi)$

Initialize ξ and ϕ randomly

while not converged **do**

Set $\underline{L}_{K_t} = 0$, ρ_t and η_t as step sizes, and $K_t \geq 0$ as a non-decreasing integer; Sample $\psi^{(k)} = T_\phi(\epsilon^{(k)})$, $\epsilon^{(k)} \sim q(\epsilon)$ for $k = 1, \dots, K_t$; take sub-sample $\mathbf{x} = \{x_i\}_{i_1:i_M}$

for $j = 1$ **to** J **do**

Sample $\psi_j = T_\phi(\epsilon_j)$, $\epsilon_j \sim q(\epsilon)$

Sample $\mathbf{z}_j = f(\tilde{\epsilon}_j, \xi, \psi_j)$, $\tilde{\epsilon}_j \sim p(\epsilon)$

$\underline{L}_{K_t} = \underline{L}_{K_t} + \frac{1}{J} \left\{ -\log \frac{1}{K_t+1} \left[\sum_{k=1}^{K_t} q_\xi(\mathbf{z}_j | \psi^{(k)}) + q_\xi(\mathbf{z}_j | \psi_j) \right] + \frac{N}{M} \log p(\mathbf{x} | \mathbf{z}_j) + \log p(\mathbf{z}_j) \right\}$

end

$t = t + 1$

$\xi = \xi + \rho_t \nabla_\xi \underline{L}_{K_t}(\{\psi^{(k)}\}_{1,K_t}, \{\psi_j\}_{1,J}, \{\mathbf{z}_j\}_{1,J})$

$\phi = \phi + \eta_t \nabla_\phi \underline{L}_{K_t}(\{\psi^{(k)}\}_{1,K_t}, \{\psi_j\}_{1,J}, \{\mathbf{z}_j\}_{1,J})$

end

A. Proofs

Proof of Inequility (3). To prove a functional form of Jensen's Inequality, let $h(\mathbf{z}) = \mathbb{E}_{\psi \sim q_\phi(\psi)} q(\mathbf{z} | \psi)$ and $\langle f, g \rangle_{L^2} = \int f(\mathbf{z})g(\mathbf{z})d\mathbf{z}$. From Theorem 1, we have convexity, and according to Theorem 6.2.1. of Kurdila & Zabaranin (2005), we have an equivalent first-order definition for convexity as

$$\text{KL}(q(\mathbf{z} | \psi) || p(\mathbf{z})) \geq \text{KL}(h(\mathbf{z}) || p) + \langle q(\mathbf{z} | \psi) - h(\mathbf{z}), \nabla_q \text{KL}(q || p)|_{h(\mathbf{z})} \rangle_{L^2}$$

Taking the expectation with respect to $\psi \sim q_\phi(\psi)$ on both sides, we have

$$\begin{aligned} & \mathbb{E}_{\psi \sim q_\phi(\psi)} \text{KL}(q(\mathbf{z} | \psi) || p(\mathbf{z})) \\ & \geq \text{KL}(h(\mathbf{z}) || p(\mathbf{z})) \\ & + \mathbb{E}_{\psi \sim q_\phi(\psi)} [\langle q(\mathbf{z} | \psi) - h(\mathbf{z}), \nabla_q \text{KL}(q || p)|_{h(\mathbf{z})} \rangle_{L^2}] \\ & = \text{KL}(h(\mathbf{z}) || p(\mathbf{z})) \\ & = \text{KL}(\mathbb{E}_{\psi \sim q_\phi(\psi)} q(\mathbf{z} | \psi) || p(\mathbf{z})). \end{aligned}$$

□

Proof of Proposition 1. We show that directly maximizing the lower bound $\underline{\mathcal{L}}$ of ELBO in (4) may drive $q(\psi)$ towards degeneracy. For VI that uses $q(\mathbf{z} | \psi)$ as its variational distribution, if supposing ψ^* is the optimum variational parameter, which means

$$\psi^* = \arg \max_{\psi} -\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{q(\mathbf{z} | \psi)}{p(\mathbf{x}, \mathbf{z})},$$

then we have

$$\begin{aligned} \underline{\mathcal{L}} & = -\mathbb{E}_{\psi \sim q_\phi(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{q(\mathbf{z} | \psi)}{p(\mathbf{x}, \mathbf{z})} \\ & = \int q_\phi(\psi) [-\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{q(\mathbf{z} | \psi)}{p(\mathbf{x}, \mathbf{z})}] d\psi \\ & \leq \int q_\phi(\psi) d\psi [-\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi^*)} \log \frac{q(\mathbf{z} | \psi^*)}{p(\mathbf{x}, \mathbf{z})}] \\ & = -\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi^*)} \log \frac{q(\mathbf{z} | \psi^*)}{p(\mathbf{x}, \mathbf{z})}. \end{aligned}$$

The equality in the above equation is reached if and only if $q(\psi) = \delta_{\psi^*}(\psi)$, which means the mixing distribution degenerates to a point mass density and hence SIVI degenerates to vanilla VI. □

Proof of Proposition 2. $B_0 = 0$ is trivial. Denote $\psi^{(0)} = \psi_v$. For iid samples $\psi^{(k)} \sim q_\phi(\psi)$, when $K \rightarrow \infty$, by the strong law of large numbers, $\tilde{h}_K(\mathbf{z}) = \frac{\sum_{k=0}^K q(\mathbf{z} | \psi^{(k)})}{K+1}$ converges almost surely to $\mathbb{E}_{q_\phi(\psi)} q(\mathbf{z} | \psi) = h_\phi(\mathbf{z})$. To prove (6), by the strong law of large numbers, we first rewrite it as the limit of a double sequence $S(K, J)$, where $K, J \in \{1, 2, \dots\}$, and check the condition for the interchange of iterated limits (Rudin, 1964; Habi, 2016): i) The double limit exists; ii) Fixing one index of the double sequence, the one side limit exists for the other index .

$$\begin{aligned} & \lim_{K \rightarrow \infty} \mathbb{E}_{\psi^{(0)}, \psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \log \frac{\sum_{k=0}^K q(\mathbf{z} | \psi^{(k)})}{K+1} \\ & = \lim_{K \rightarrow \infty} \lim_{J \rightarrow \infty} \frac{1}{J} \sum_{j=1}^J \log \frac{1}{K+1} \sum_{k=0}^K q(\mathbf{z} | \psi_j^{(k)}) \\ & \triangleq \lim_{K \rightarrow \infty} \lim_{J \rightarrow \infty} S(K, J). \end{aligned}$$

Here $\psi_j^{(k)}$ are iid samples from $q(\psi)$. For i) we show double limit $\lim_{K, J \rightarrow \infty} S(K, J) = \log h(\mathbf{z})$. For $\forall \epsilon > 0$, $\exists N(\epsilon)$, when $K, J > N(\epsilon)$, $|\log \frac{1}{K+1} \sum_{k=0}^K q(\mathbf{z} | \psi_j^{(k)}) - \log h(\mathbf{z})| < \epsilon$ thanks to the law of large numbers, then

$$\begin{aligned} & \left| \sum_{j=1}^J \log \frac{1}{K+1} \sum_{k=0}^K q(\mathbf{z} | \psi_j^{(k)}) - J \log h(\mathbf{z}) \right| \\ & \leq \sum_{j=1}^J \left| \log \frac{1}{K+1} \sum_{k=0}^K q(\mathbf{z} | \psi_j^{(k)}) - \log h(\mathbf{z}) \right| \leq J\epsilon. \end{aligned}$$

Deviding both sides by J we get $|S(K, J) - \log h(\mathbf{z})| \leq \epsilon$ when $K, J > N(\epsilon)$. By definition, we have $\lim_{K, J \rightarrow \infty} S(K, J) = \log h(\mathbf{z})$.
 ii) for each fixed $J \in \mathbb{N}$, $\lim_{K \rightarrow \infty} S(K, J) = \log h(\mathbf{z})$ exists; for each fixed $K \in \mathbb{N}$, $\lim_{J \rightarrow \infty} S(K, J) = \mathbb{E}_{\psi^{(0)}, \psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \log \frac{\sum_{k=0}^K q(\mathbf{z} | \psi^{(k)})}{K+1} \leq \log h(\mathbf{z})$ also exists. The limitation can then be interchanged and (6) is proved. Therefore, we have

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathcal{L}_k &= \underline{\mathcal{L}} + \mathbb{E}_{\psi} \text{KL}(q(\mathbf{z} | \psi) || h_{\phi}(\mathbf{z})) \\ &= \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \left[\log \frac{q(\mathbf{z} | \psi)}{h_{\phi}(\mathbf{z})} - \log \frac{q(\mathbf{z} | \psi)}{p(x, \mathbf{z})} \right] \\ &= -\mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \psi)} \log \frac{h_{\phi}(\mathbf{z})}{p(x, \mathbf{z})} \\ &= -\mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \log \frac{h_{\phi}(\mathbf{z})}{p(x, \mathbf{z})} = \mathcal{L}. \end{aligned}$$

□

Proof of Proposition 3. Assume integer $K > M > 0$. Let \mathcal{I} be the set that consists of all the subsets of $\{1, \dots, K\}$ with cardinality M . Let I be a discrete uniform random variable and for element $\{i_1, \dots, i_M\} \in \mathcal{I}$, $P(I = \{i_1, \dots, i_M\}) = \frac{1}{\binom{K}{M}}$. We have $\mathbb{E}_I \frac{1}{M} \sum_{i \in I} q(\mathbf{z} | \psi^i) = \frac{1}{K} \sum_{i=1}^K q(\mathbf{z} | \psi^i)$. To show $\bar{\mathcal{L}}_K = \bar{\mathcal{L}} - A_K$ is monotonic decreasing, we only need to show A_K is monotonic increasing:

$$\begin{aligned} A_K &= \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \\ &\quad \log \frac{\frac{1}{K} \sum_{i=1}^K q(\mathbf{z} | \psi^{(i)})}{q(\mathbf{z} | \psi)} \\ &= \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \\ &\quad \log \mathbb{E}_I \left[\frac{\frac{1}{M} \sum_{i \in I} q(\mathbf{z} | \psi^{(i)})}{q(\mathbf{z} | \psi)} \right] \\ &\geq \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \\ &\quad \mathbb{E}_I \log \frac{\frac{1}{M} \sum_{i \in I} q(\mathbf{z} | \psi^{(i)})}{q(\mathbf{z} | \psi)} \\ &= \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(M)} \sim q(\psi)} \\ &\quad \log \frac{\frac{1}{M} \sum_{i=1}^M q(\mathbf{z} | \psi^{(i)})}{q(\mathbf{z} | \psi)} \\ &= A_M. \end{aligned}$$

We now show $\lim_{K \rightarrow \infty} \bar{\mathcal{L}}_K = \mathcal{L}$. Again, by the strong law of large numbers, $\frac{1}{K} \sum_{i=1}^K q(\mathbf{z} | \psi^{(i)})$ converges almost

surely to $\mathbb{E}_{\psi \sim q(\psi)} q(\mathbf{z} | \psi) = h_{\phi}(\mathbf{z})$ and hence

$$\begin{aligned} \lim_{K \rightarrow \infty} \bar{\mathcal{L}}_K &= \bar{\mathcal{L}} + \mathbb{E}_{\psi} \text{KL}(h_{\phi}(\mathbf{z}) || q(\mathbf{z} | \psi)) \\ &= -\mathbb{E}_{\mathbf{z} \sim h_{\phi}(\mathbf{z})} \mathbb{E}_{\psi \sim q(\psi)} \left[\log \frac{q(\mathbf{z} | \psi)}{p(x, \mathbf{z})} + \log \frac{h_{\phi}(\mathbf{z})}{q(\mathbf{z} | \psi)} \right] \\ &= \mathcal{L}. \end{aligned}$$

□

Proof of Equation (11). The gradient of B_K with respect to ϕ can be expressed as

$$\begin{aligned} \nabla_{\phi} B_K &= \nabla_{\phi} \mathbb{E}_{\psi \sim q(\psi)} \mathbb{E}_{\psi^{(1)}, \dots, \psi^{(K)} \sim q(\psi)} \left[\text{KL} \left(q(\mathbf{z} | \psi) \left\| \frac{q(\mathbf{z} | \psi) + \sum_{k=1}^K q(\mathbf{z} | \psi^{(k)})}{K+1} \right. \right) \right] \\ &= \mathbb{E}_{\epsilon, \epsilon^{(1)}, \dots, \epsilon^{(K)} \sim p(\epsilon)} \nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | T_{\phi}(\epsilon))} \left[\log \frac{q(\mathbf{z} | T_{\phi}(\epsilon))}{\frac{q(\mathbf{z} | T_{\phi}(\epsilon)) + \sum_{k=1}^K q(\mathbf{z} | T_{\phi}(\epsilon^{(k)}))}{K+1}} \right] \\ &= \mathbb{E}_{\epsilon, \dots, \epsilon^{(K)}} \nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | T_{\phi}(\epsilon))} \log \left[r_{\phi}(\mathbf{z}, \epsilon, \epsilon^{(1:K)}) \right] \\ &= \mathbb{E}_{\epsilon, \dots, \epsilon^{(K)} \sim p(\epsilon)} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | T_{\phi}(\epsilon))} \left\{ q(\mathbf{z} | T_{\phi}(\epsilon)) \nabla_{\phi} \log \left[r_{\phi}(\mathbf{z}, \epsilon, \epsilon^{(1:K)}) \right] \right. \\ &\quad \left. + \left[\nabla_{\phi} \log q(\mathbf{z} | T_{\phi}(\epsilon)) \right] \log \left[r_{\phi}(\mathbf{z}, \epsilon, \epsilon^{(1:K)}) \right] \right\}. \end{aligned}$$

□

B. Bayesian Logistic Regression

We consider datasets *waveform* ($n = 5000$, $V = 21$, and $400/4600$ for training/testing), *spam* ($n = 3000$, $V = 2$, and $2000/1000$ for training/testing), and *nodal* ($n = 53$, $V = 5$, and $25/28$ for training/testing). The training-set-size to feature-dimension ratio n_{train}/V varies in these three datasets, and we expect the posterior uncertainty to be large if this ratio is small.

The contribution of observation i to the likelihood can be expressed as

$$P(y_i | \mathbf{x}_i, \beta) = \frac{e^{y \mathbf{x}'_i \beta}}{1 + e^{\mathbf{x}'_i \beta}} \propto e^{(y - \frac{1}{2}) \mathbf{x}'_i \beta} \mathbb{E}_{\omega_i} \left[e^{-\frac{\omega_i (\mathbf{x}'_i \beta)^2}{2}} \right],$$

where the expectation is taken respect to a Pólya-Gamma (PG) distribution (Polson et al., 2013) as $\omega_i \sim \text{PG}(1, 0)$, and hence we have an augmented likelihood as

$$P(y_i, \omega_i | \mathbf{x}_i, \beta) \propto e^{(y - \frac{1}{2}) \mathbf{x}'_i \beta - \frac{1}{2} \omega_i (\mathbf{x}'_i \beta)^2}.$$

B.1. Gibbs Sampling via Data Augmentation

Denoting $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)'$, $\mathbf{y} = (y_1, \dots, y_N)'$, $\mathbf{A} = \text{diag}(\alpha_0, \dots, \alpha_V)'$, and $\mathbf{\Omega} = \text{diag}(\omega_1, \dots, \omega_N)$, we have

$$(\omega_i | -) \sim \text{PG}(1, \mathbf{x}'_i \boldsymbol{\beta}), \quad (\boldsymbol{\beta} | -) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma} = (\mathbf{A} + \mathbf{X}'\mathbf{\Omega}\mathbf{X})^{-1}$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma}\mathbf{X}'(\mathbf{y} - 1/2)$. To sample from the Pólya-Gamma distribution, a random sample from which can be generated as a weighted sum of an infinite number of *iid* gamma random variables, we follow Zhou (2016) to truncate the infinite sum to the summation of M gamma random variables, where the parameters of the M th gamma random variable are adjusted to match the mean and variance of the finite sum with those of the infinite sum. We set $M = 5$ in this paper.

B.2. Mean-Field Variational Inference with Diagonal Covariance Matrix

We choose a fully factorized Q distribution as

$$Q = \left[\prod_i q(\omega_i) \right] \left[\prod_v q(\beta_v) \right].$$

To exploit conjugacy, defining

$$\begin{aligned} q(\omega_i) &= \text{PG}(1, \lambda_i), \\ q(\beta_v) &= \mathcal{N}(\mu_v, \sigma_v^2), \end{aligned}$$

we have closed-form coordinate ascent variational inference update equations as

$$\begin{aligned} \lambda_i &= \sqrt{\mathbb{E}[(\mathbf{x}'_i \boldsymbol{\beta})^2]} = \sqrt{\mathbf{x}'_i \mathbb{E}[\boldsymbol{\beta} \boldsymbol{\beta}'] \mathbf{x}_i}, \\ \sigma_v^2 &= \left(\mathbb{E}[\alpha_v] + \sum_i \mathbb{E}[\omega_i] x_{iv}^2 \right)^{-1} \\ \mu_v &= \sigma_v^2 \sum_i x_{iv} \left\{ y_i - 1/2 - \mathbb{E}[\omega_i] \sum_{\tilde{v} \neq v} x_{i\tilde{v}} \mathbb{E}[\beta_{\tilde{v}}] \right\}, \end{aligned}$$

where the expectations with respect to the q distributions can be expressed as $\mathbb{E}[\boldsymbol{\beta} \boldsymbol{\beta}'] = \boldsymbol{\mu} \boldsymbol{\mu}' + \text{diag}(\sigma_0^2, \dots, \sigma_V^2)$ and $\mathbb{E}[\omega_i] = \tanh(\lambda_i/2)/(2\lambda_i)$.

B.3. Mean-Field Variational Inference with Full Covariance Matrix

We choose a fully factorized Q distribution as

$$Q = \left[\prod_i q(\omega_i) \right] q(\boldsymbol{\beta}).$$

To exploit conjugacy, defining

$$\begin{aligned} q(\omega_i) &= \text{PG}(1, \lambda_i), \\ q(\boldsymbol{\beta}) &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \end{aligned}$$

we have closed-form coordinate ascent variational inference update equations as

$$\begin{aligned} \lambda_i &= \sqrt{\mathbb{E}[(\mathbf{x}'_i \boldsymbol{\beta})^2]} = \sqrt{\mathbf{x}'_i \mathbb{E}[\boldsymbol{\beta} \boldsymbol{\beta}'] \mathbf{x}_i}, \\ \boldsymbol{\Sigma} &= (\mathbb{E}[\mathbf{A}] + \mathbf{X}' \mathbb{E}[\boldsymbol{\Omega}] \mathbf{X})^{-1}, \quad \boldsymbol{\mu} = \boldsymbol{\Sigma} \mathbf{X}' (\mathbf{y} - 1/2), \end{aligned}$$

where the expectations with respect to the q distributions can be expressed as $\mathbb{E}[\boldsymbol{\beta} \boldsymbol{\beta}'] = \boldsymbol{\mu} \boldsymbol{\mu}' + \boldsymbol{\Sigma}$ and $\mathbb{E}[\omega_i] = \tanh(\lambda_i/2)/(2\lambda_i)$. Note the update equations shown above are identical to those shown in Jaakkola & Jordan (2000).

B.4. SIVI Configuration

For inputs in Algorithm 1, we choose a multi-layer perceptron with layer size [100, 200, 100] as T_ϕ for $\boldsymbol{\psi} = T_\phi(\boldsymbol{\epsilon})$, $\boldsymbol{\epsilon}$ as 50 dimensional isotropic Gaussian random variable and $K = 100$, $J = 50$. For the explicit layer, we choose an MVN as $q_\xi(\mathbf{z} | \boldsymbol{\psi}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\psi}, \boldsymbol{\xi})$. In this setting, $\boldsymbol{\psi}$ is the mean variable mixed with implicit distribution $q_\phi(\boldsymbol{\psi})$ while $\boldsymbol{\xi}$ is the covariance matrix which can be either diagonal or full. In the experiments, we update the neural network parameter ϕ by the Adam optimizer, with learning rate 0.01. We update $\boldsymbol{\xi}$ by gradient ascent, with step size $\eta_t = 0.001 * 0.9^{\text{iteration}/100}$. The implicit layer parameter ϕ and explicit layer parameter $\boldsymbol{\xi}$ are updated iteratively.

C. Experimental Settings and Results for SIVAE

We implement SIVI with $M = 3$ stochastic hidden layers, with the dimensions of hidden layers $[\ell_1, \ell_2, \ell_3]$ as [150, 150, 150] and with the dimensions of injected noises $[\epsilon_1, \epsilon_2, \epsilon_3]$ as [150, 100, 50]. Between two adjacent stochastic layers there is a fully connected deterministic layer with size 500 and *ReLU* activation function. We choose binary pepper and salt noise (Im et al., 2017) for $q_t(\boldsymbol{\epsilon})$. The model is trained for 2000 epochs with mini-batch size 200 and step-size $0.001 * 0.75^{\text{epoch}/100}$. K_t is gradually increased from 1 to 100 during the first 1500 epochs. The explicit and implicit layers are trained iteratively. Warm-up is used during the first 300 epochs as suggested by Sønderby et al. (2016) to gradually impose the prior regularization term $\text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$. The model is trained end-to-end using the Adam optimizer. After training process, as in Rezende et al. (2014) and Burda et al. (2015), we compute the marginal likelihood for test set by importance sampling with $S = 2000$:

$$\log p(\mathbf{x}) \approx \log \frac{1}{S} \sum_{s=1}^S \frac{p(\mathbf{x} | \mathbf{z}_s) p(\mathbf{z}_s)}{\hat{h}(\mathbf{z}_s | \mathbf{x})}, \quad \mathbf{z}_s \sim h(\mathbf{z}_s | \mathbf{x}),$$

where

$$\hat{h}(\mathbf{z}_s | \mathbf{x}) = \frac{1}{M} \sum_{k=1}^M q(\mathbf{z}_s | \boldsymbol{\psi}^{(k)}), \quad \boldsymbol{\psi}^{(k)} \stackrel{iid}{\sim} q_\phi(\boldsymbol{\psi} | \mathbf{x})$$

is used to estimate $h(\mathbf{z}_s | \mathbf{x})$; we set $M = 100$. The performance of SIVI and a comparison to reported results with other methods are provided in Table 2.

Table 2. Comparison of the negative log evidence between various algorithms.

Methods	$-\log p(\mathbf{x})$
<i>Results below form Burda et al. (2015)</i>	
VAE + IWAE	= 86.76
IWAE + IWAE	= 84.78
<i>Results below form Salimans et al. (2015)</i>	
DLGM + HVI (1 leapfrog step)	= 88.08
DLGM + HVI (4 leapfrog step)	= 86.40
DLGM + HVI (8 leapfrog steps)	= 85.1
<i>Results below form Rezende & Mohamed (2015)</i>	
DLGM+NICE (Dinh et al., 2014) (k = 80)	≤ 87.2
DLGM+NF (k = 40)	≤ 85.7
DLGM+NF (k = 80)	≤ 85.1
<i>Results below form Gregor et al. (2015)</i>	
DLGM	≈ 86.60
NADE	= 88.33
DBM 2hl	≈ 84.62
DBN 2hl	≈ 84.55
EoNADE-5 2hl (128 orderings)	= 84.68
DARN 1hl	≈ 84.13
<i>Results below form Maaløe et al. (2016)</i>	
Auxiliary VAE (L=1, IW=1)	≤ 84.59
<i>Results below form Mescheder et al. (2017)</i>	
VAE + IAF (Kingma et al., 2016)	≈ 84.9 ± 0.3
Auxiliary VAE (Maaløe et al., 2016)	≈ 83.8 ± 0.3
AVB + AC	≈ 83.7 ± 0.3
SIVI (3 stochastic layers)	= 84.07
SIVI (3 stochastic layers)+ IW($\tilde{K} = 10$)	= 83.25

D. Additional Figures

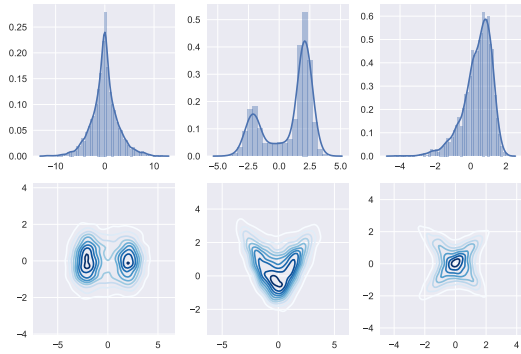


Figure 8. Visualization of the MLP based implicit distributions $\psi \sim q(\psi)$, which are mixed with isotropic Gaussian (or Log-Normal) distributions to approximate the target distributions shown in Figure 1.

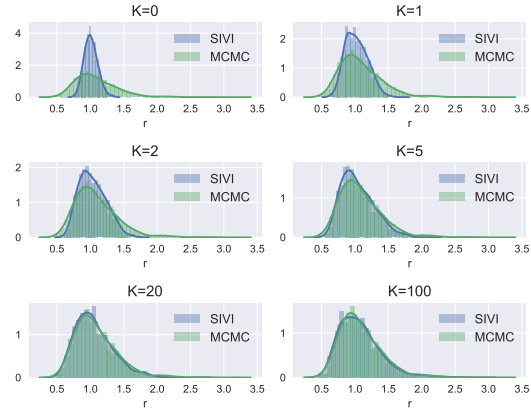


Figure 9. The marginal posterior distribution of the negative binomial dispersion parameter r inferred by SIVI becomes more accurate as K increases

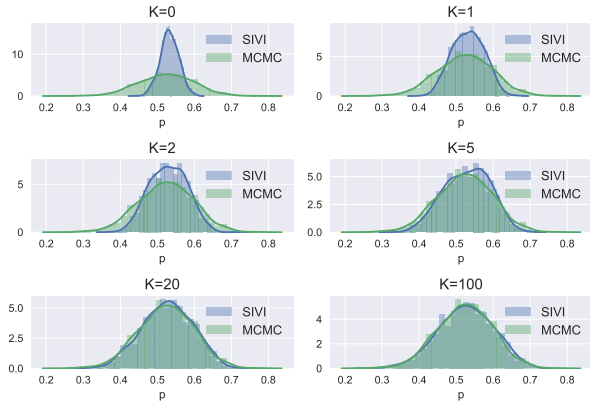


Figure 10. The marginal posterior distribution of the negative binomial probability parameter p inferred by SIVI becomes more accurate as K increases.

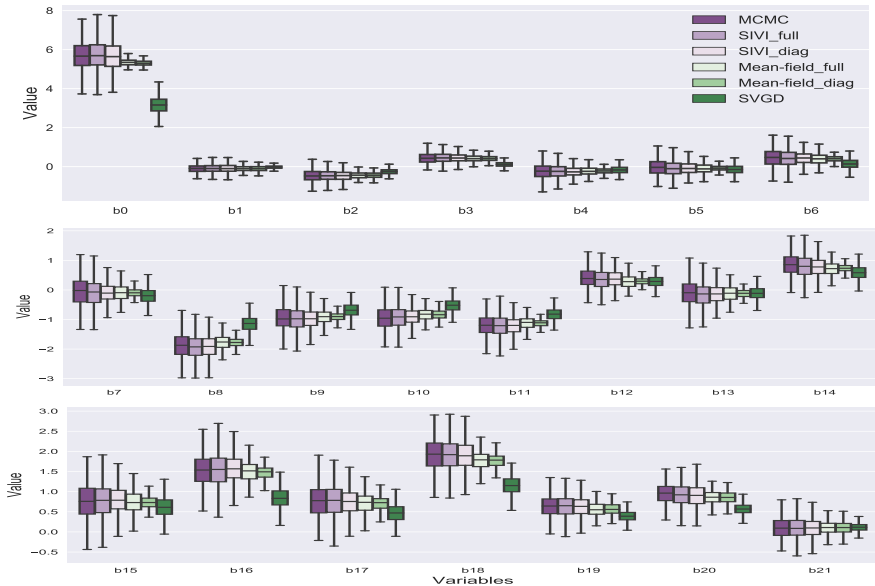


Figure 11. Comparison of all marginal posteriors of β_v inferred by various methods for Bayesian logistic regression on *waveform*.

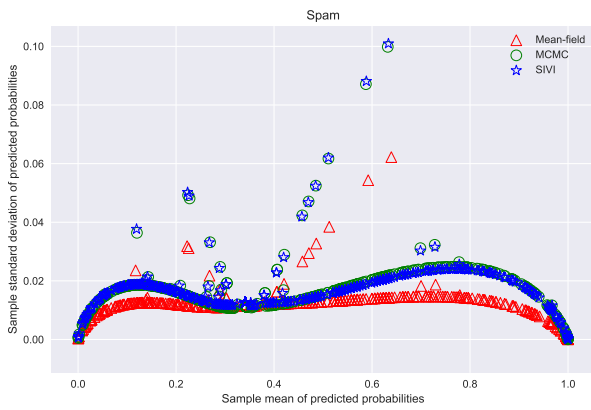


Figure 12. Sample means and standard deviations of predictive probabilities for dataset *spam*.

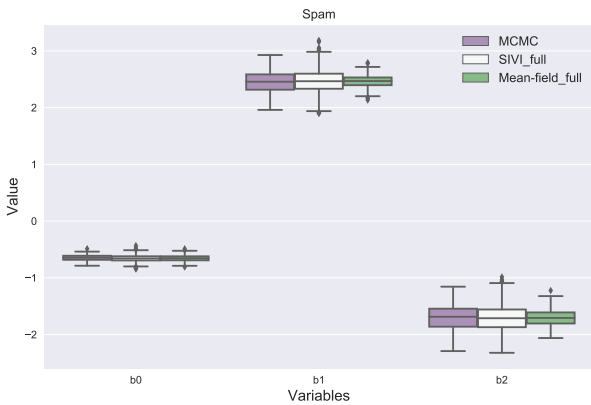


Figure 13. Boxplot of marginal posteriors inferred by MCMC, SIVI, and MFVI for dataset *spam*.

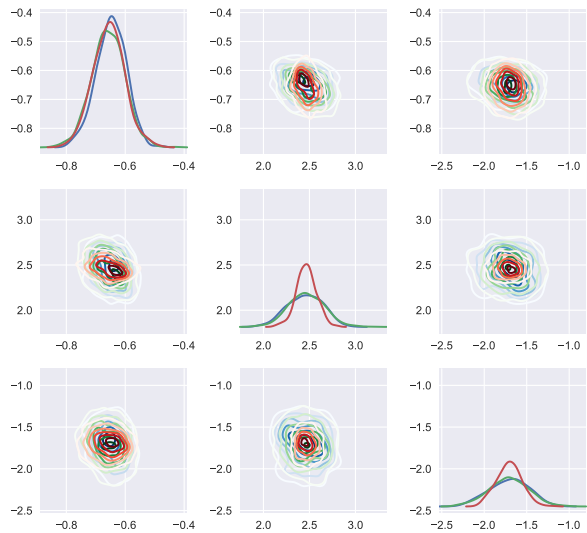


Figure 14. Univariate marginal and pairwise joint posteriors for dataset *spam*. Blue, green, and red are for MCMC, SIVI with a full covariance matrix, and MFVI with a full covariance matrix, respectively.

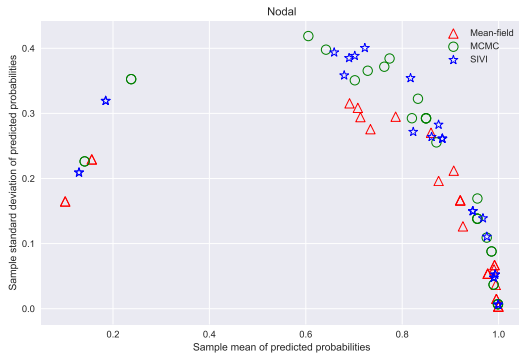


Figure 15. Sample means and standard deviations of predictive probabilities for dataset *nodal*.

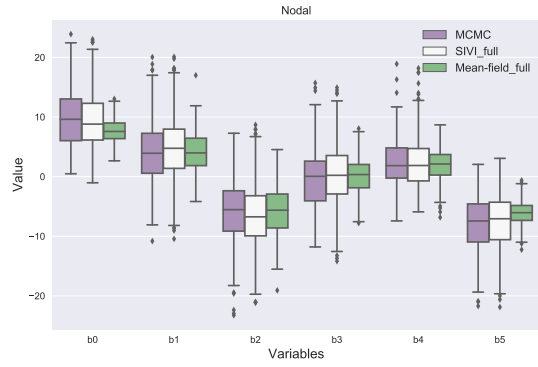


Figure 16. Boxplot of marginal posteriors inferred by MCMC, SIVI, and MFVI for dataset *nodal*.

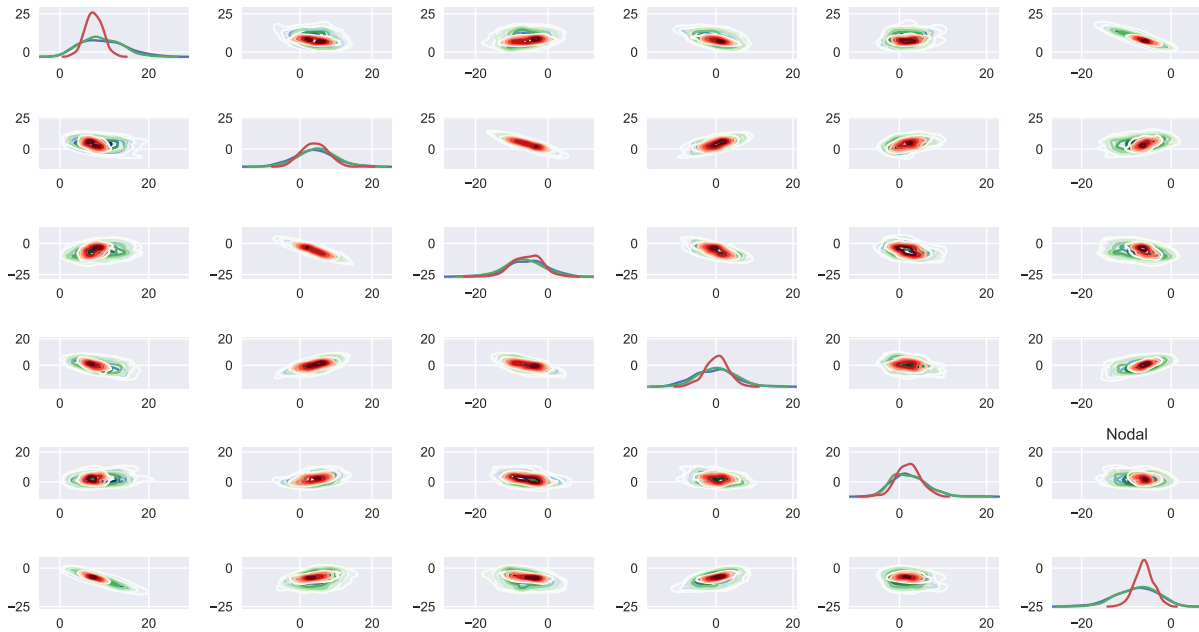


Figure 17. Univariate marginal and pairwise joint posteriors for dataset *nodal*. Blue, green, and red are for MCMC, SIVI with a full covariance matrix, and MFVI with a full covariance matrix, respectively.